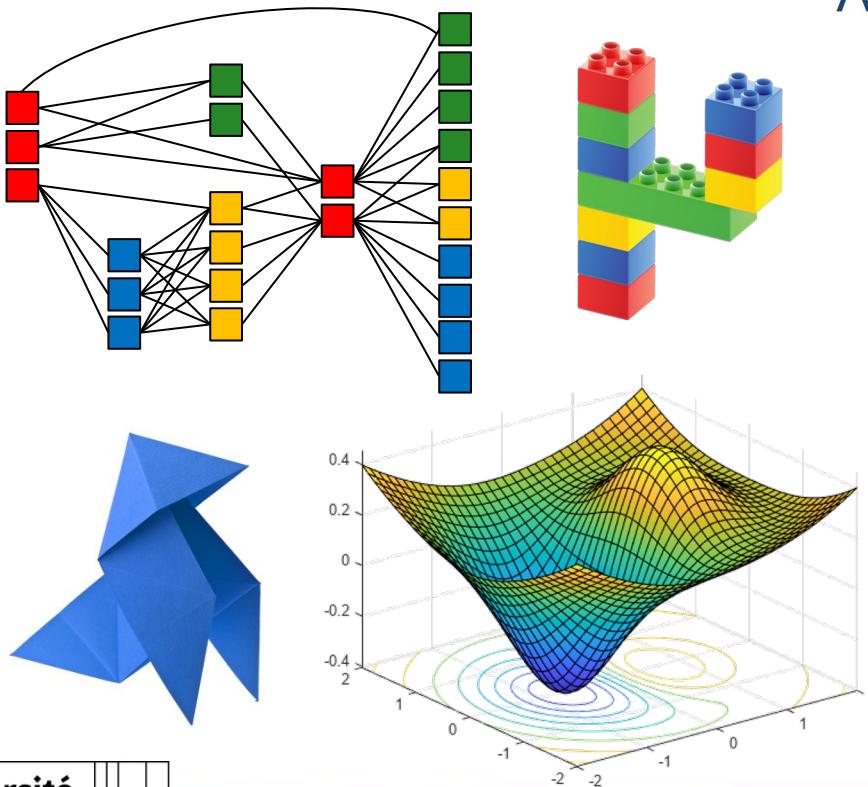


# Autumn School Series in Acoustics

## Lecture 1: Fundamentals of Machine Learning

Eindhoven University of Technology  
3-7 November 2025

Antoine Deleforge & Mathieu Gaborit



# Outline

---

## 1. Introduction

AI? Machine Learning? Neural Networks?

## 2. Machine Learning

Definition, Characteristics, Types, Algorithms

## 3. Deep Neural Networks

Definition, Motivation, How to Train Them

## 4. Limits of Machine Learning

Overfitting, Data Biases

Then: Tutorial to put this in practice!

## Artificial Intelligence

Machine Learning

Neural Networks

Generative Models

Deep Learning



*ChatGPT, Grok, MidJourney,  
Stable Diffusion, ...*

## Artificial Intelligence

- What is intelligence? (Psychology / Philosophy)
- AI roots: collective imagination (Science fiction)
- Used for communication / marketing / fundraising
- AI ethics and safety (law, philosophy, politics)

### “Applied” A.I.

- Solving « difficult » numerical problems (*optimization, game theory, ...*)

- Representing knowledge and reasoning (logic, planning, inference, ...)
- Managing uncertainty (representation, probabilities, decision procedures, ...)
- Language and communication (natural language processing, text ↔ speech, avatars, robotics, ...)

#### Symbolic A.I.

- Explicit knowledge representation
- Rule-based methods
- *Deep Blue* chess engine (1997)
- High-level programming languages (Python, mathematica)

#### Machine Learning

- **Algorithms** that (implicitly) build **models** from **data** to achieve **tasks**

#### Neural Networks

- Class of **parameterized** families of **non-linear** functions

#### Deep Learning

- « Deep » neural networks applied to machine learning

#### Generative Models

- The model is a **probability distribution** that can be **sampling**

PINNs  
...

“A.I.” / AI model (since ~2023)

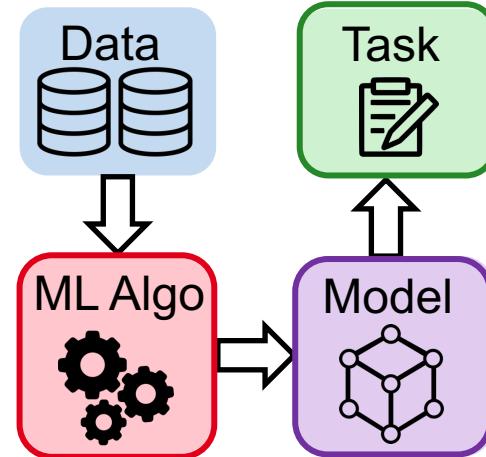
*ChatGPT, Grok, MidJourney,  
Stable Diffusion, ...*

### Definition

The study of **Algorithms** that build **Models** from **Data** to achieve **Tasks**.

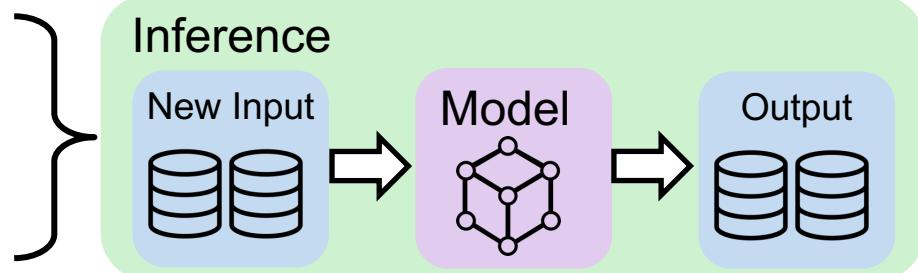
#### Model (in machine learning):

- A “*program created by a program*”
- Can be **deterministic** or **stochastic**
- **Data-Driven** as opposed to **Physics- / Knowledge-Driven** model



#### Tasks:

- Sort / Visualize / Represent the **Data**
- **Infer** from **new input Data**:
  - Estimate explanatory quantities
  - Generate, complete, predict
  - Transform, convert, translate
  - Decide

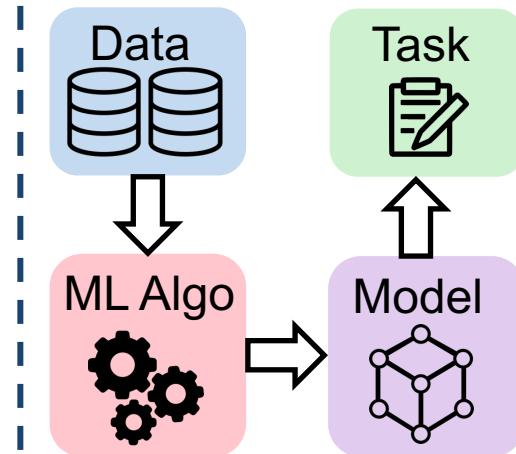
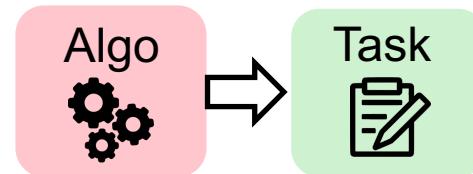


Not considered ML,  
even if the same Task  
is performed!

## Contrast it with *symbolic A.I.*

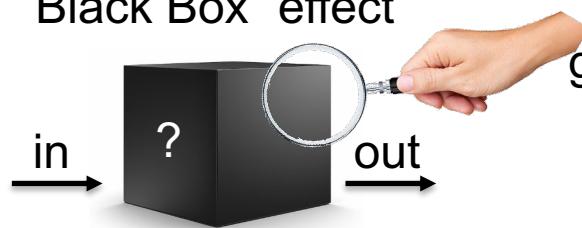
```
def add_numbers(a, b):
    return a + b

def square(num):
    result = num * num
    return result
```



- ML models do not necessarily contain **explicit** concepts with **meanings**

“Black Box” effect



Explainability / Interpretability / Theoretical guarantees for models = a rich subfield of ML  
=> but done ***a posteriori***

- Is Machine Learning **enough** for A.I.?



- A heated debate! Ex: Yann LeCun vs. Gary Marcus (NYU)
- **Neuro-symbolic A.I.** aims at combining Symbolic A.I. and Neural Networks
- Attractive approach: Hybridizing **physics-** and **data**-driven models

# Types of ML based on available data

- Supervised Learning

Ex 1: Sound Source Localization

- A dataset of audio recordings, annotated by source position
- Given a new recording, where's the sound?

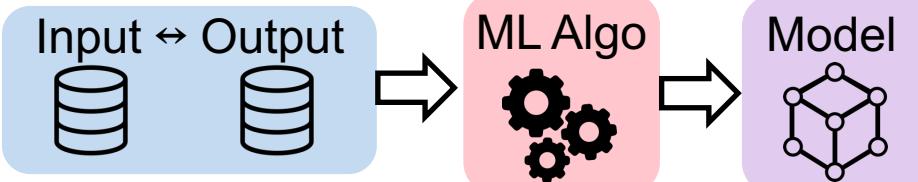


Ex 2: Sound Event Recognition

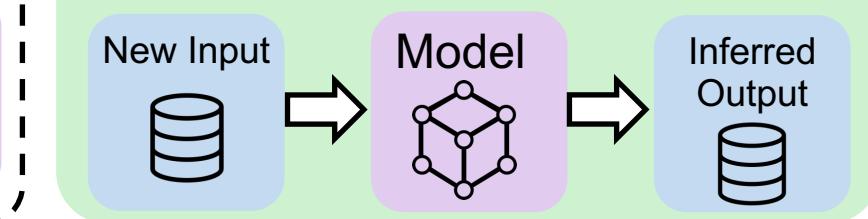
- An annotated dataset of cat sounds and baby sounds
- Given a new recording, is it a cat or a baby?



## Training Stage



## Inference / Evaluation Stage

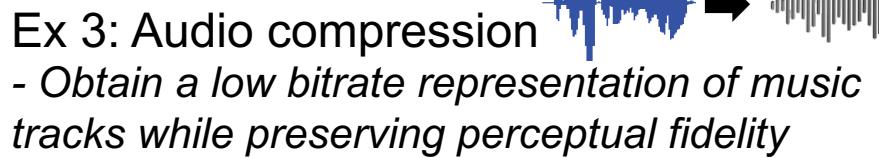
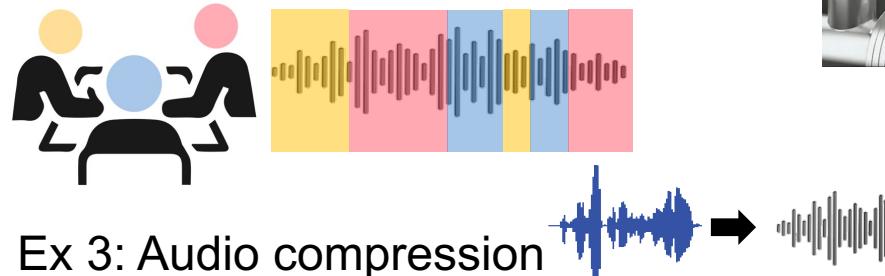


# Types of ML based on available data

- Unsupervised Learning

## Ex 1: Speaker diarization

- Given a 1h audio recording of a meeting between 3 unknown people, annotate the speaking turns



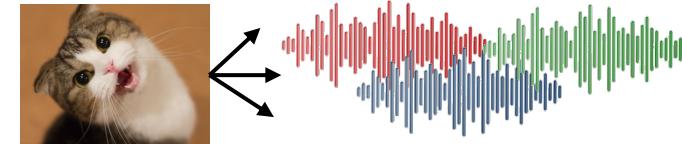
## Ex 2: Anomaly detection

- Identify a gas leak, or a train rail defect from acoustic recordings (industry monitoring)



## Ex 4: Audio generation

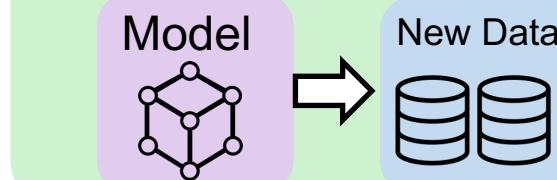
- Given a dataset of cat recordings, generate new plausible cat sounds



### Training Stage

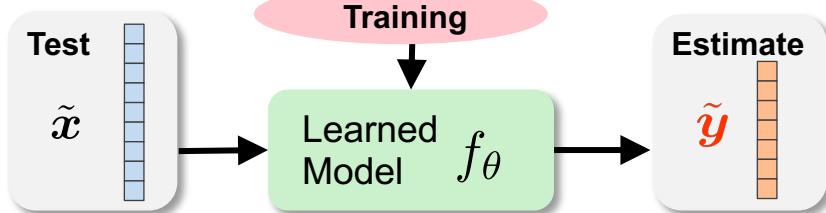
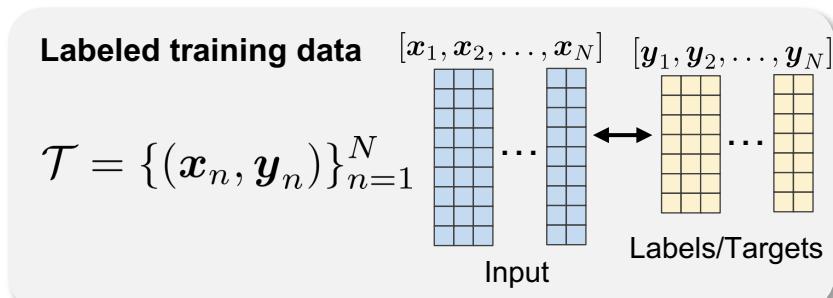


### Inference / Evaluation Stage

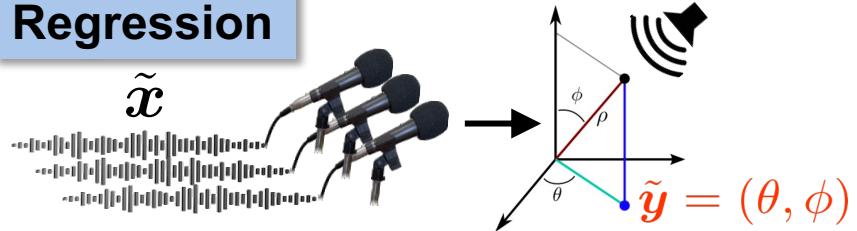


## Supervised Learning

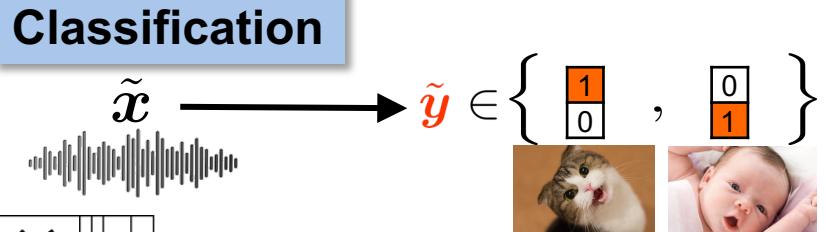
OUTPUT



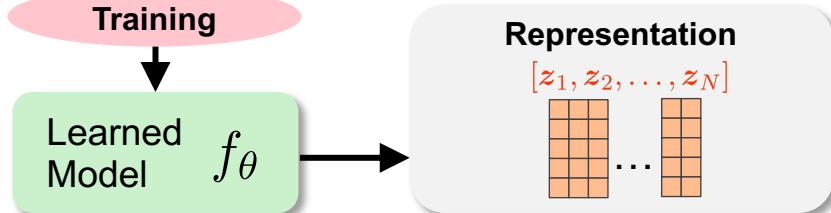
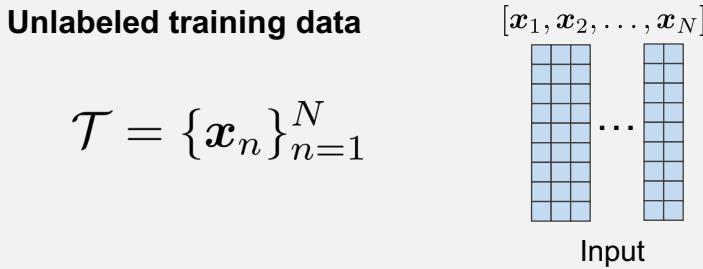
CONTINUOUS



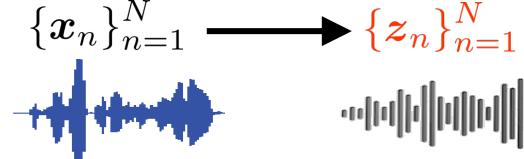
DISCRETE



## Unsupervised (representation) Learning



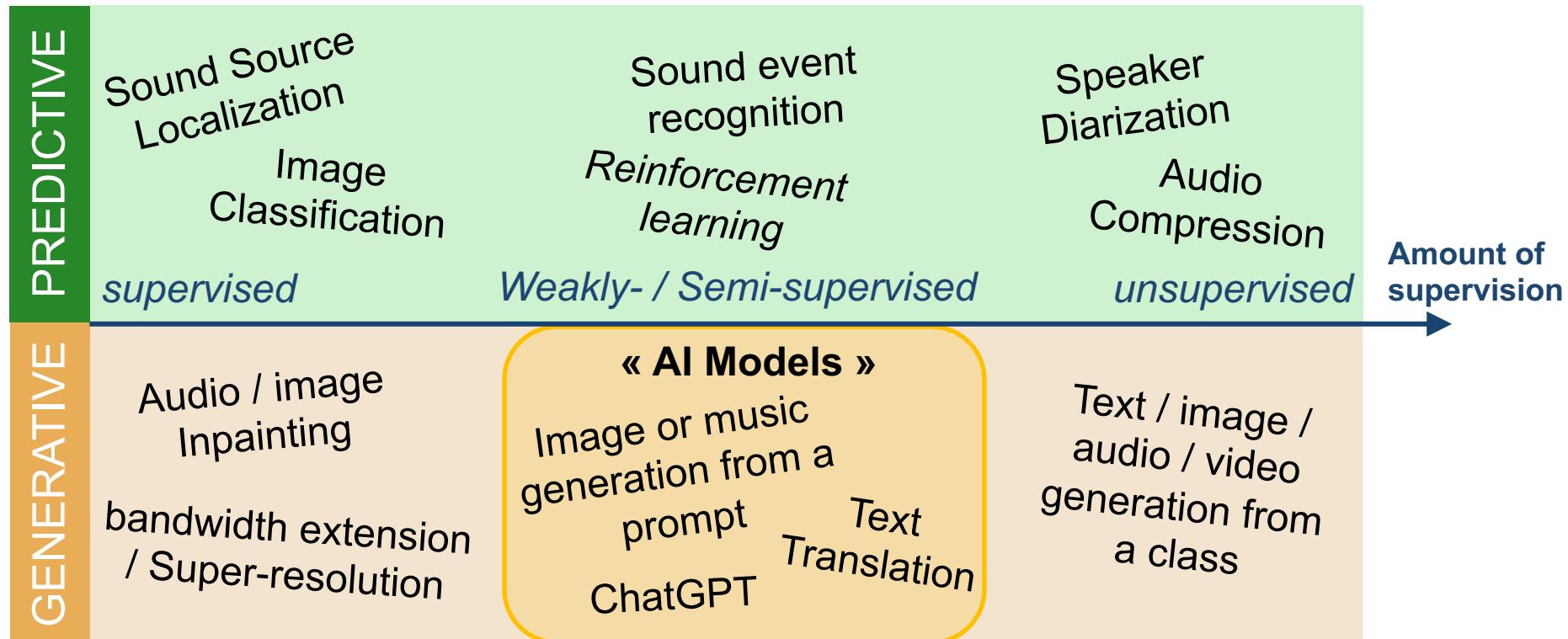
## Dimensionality Reduction



## Clustering

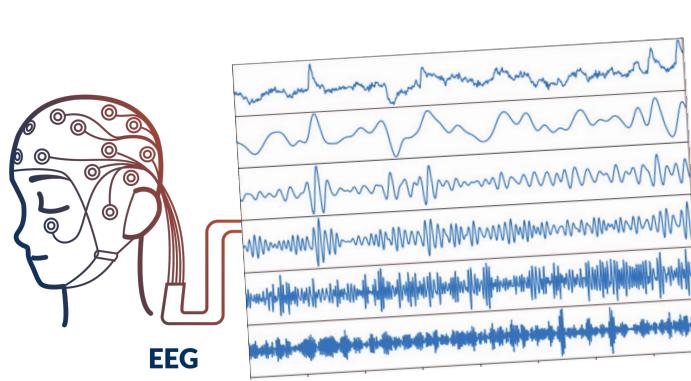
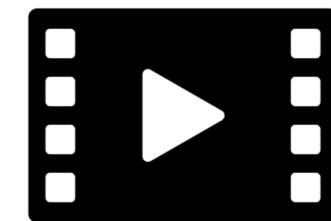
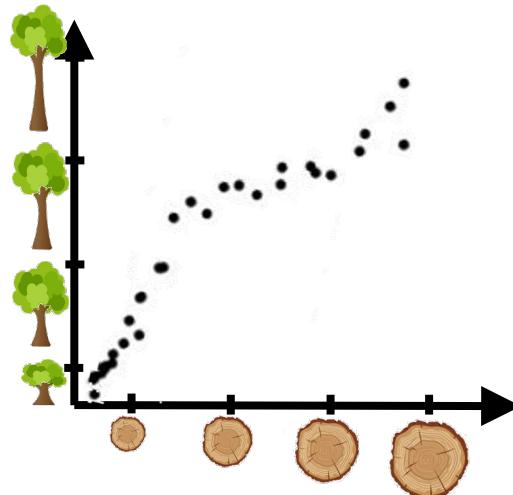


# A “Continuum” of Supervision



# Types of data used in machine learning

- Numerical Data / Continuous Data / Signals / Vectors in  $\mathbb{R}^D$



# Types of data used in machine learning

- Categorical Data

ImageNet Labels

airplane  
automobile  
bird  
cat  
deer  
dog  
frog  
horse  
ship  
truck  
...

Sex  $\in \{\text{male, female, other}\}$   
ZIP  $\in \{67000, 75009, \dots\}$   
Country  $\in \{\text{France, Spain, ...}\}$

- Heterogenous / Tabular Data



- Text Data

Artificial intelligence

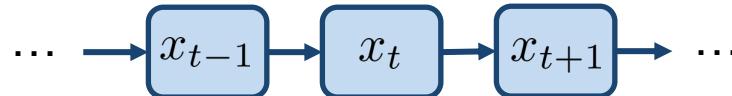
Article Talk

From Wikipedia, the free encyclopedia

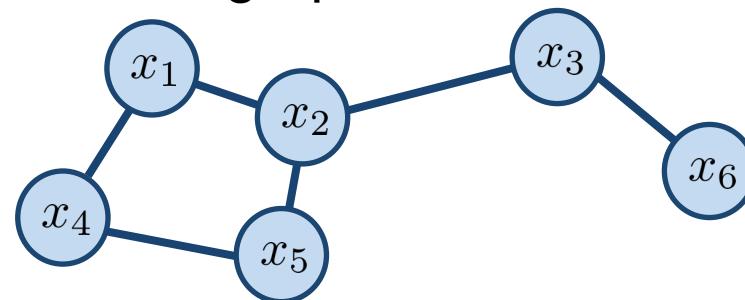
**Artificial intelligence** (AI) is **intelligence**—perceiving, synthesizing, and inferring information—demonstrated by **machines**, as opposed to intelligence displayed by **non-human animals** and **humans**. Example tasks in which this is done include speech recognition, computer vision, translation between (natural) languages, as well as other mappings of inputs.

AI applications include advanced **web search engines** (e.g., **Google Search**), **recommendation systems** (used by **YouTube**, **Amazon** and **Netflix**), **understanding human speech** (such as **Siri** and **Alexa**), **self-driving cars** (e.g., **Waymo**), **generative or creative tools** (**ChatGPT** and **AI art**), **automated decision-making** and competing at the highest level in **strategic game systems** (such as **chess** and **Go**).<sup>[1]</sup>

- Time series

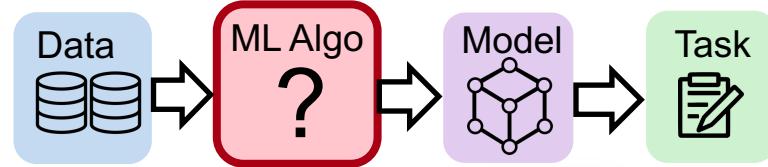


- Data on graphs



# Machine learning algorithms

*How to make a program that makes a program?*



- Search in the set of all python functions?

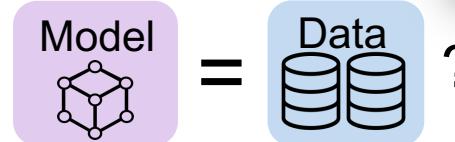
- a gigantic combinatorial set

- Simply **memorize** the data

- “Lazy Learning”

- Ex: k-NN, look-up table, naive Bayes, case-based

- Slow, large storage, non-robust



```

def add_numbers(a, b):
    return a + b
def square(num):
    result = num * num
    return result
  
```

## Model Fitting:

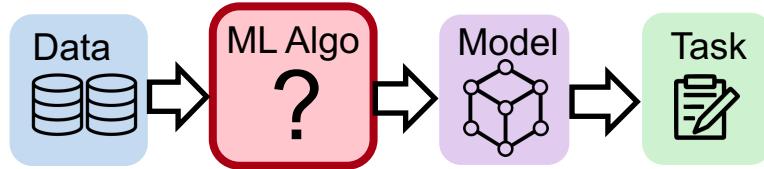
Find the *best* model within a **parameterized family**  $\mathcal{F} = \{m_\theta\}_{\theta \in \Theta}$



- $m_\theta$  = a jean
- $\theta$  = its (width, length)
- $\mathcal{F}$  = the shelves

# Machine learning algorithms

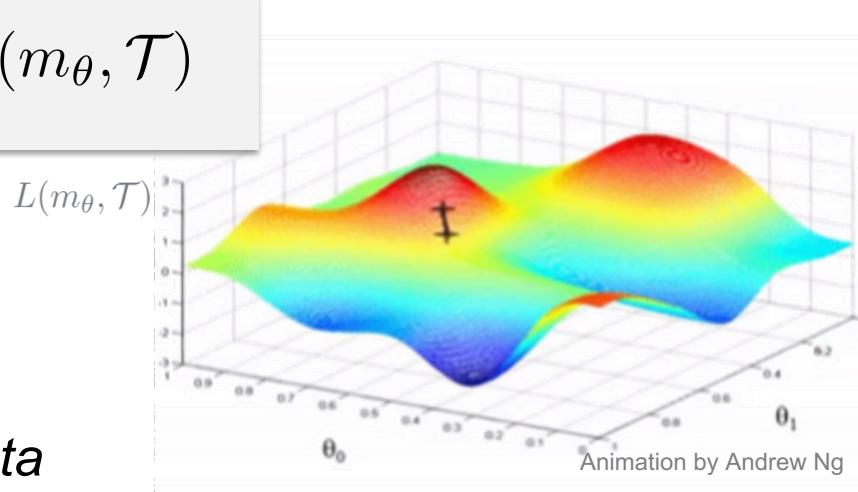
*How to make a program that makes a program?*



- The model is found by minimizing a **total loss / cost function**  $L$  over the set of parameters, for a given training **dataset**  $\mathcal{T}$ :

$$\hat{m} = m_{\hat{\theta}} \quad \text{where} \quad \hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} L(m_\theta, \mathcal{T})$$

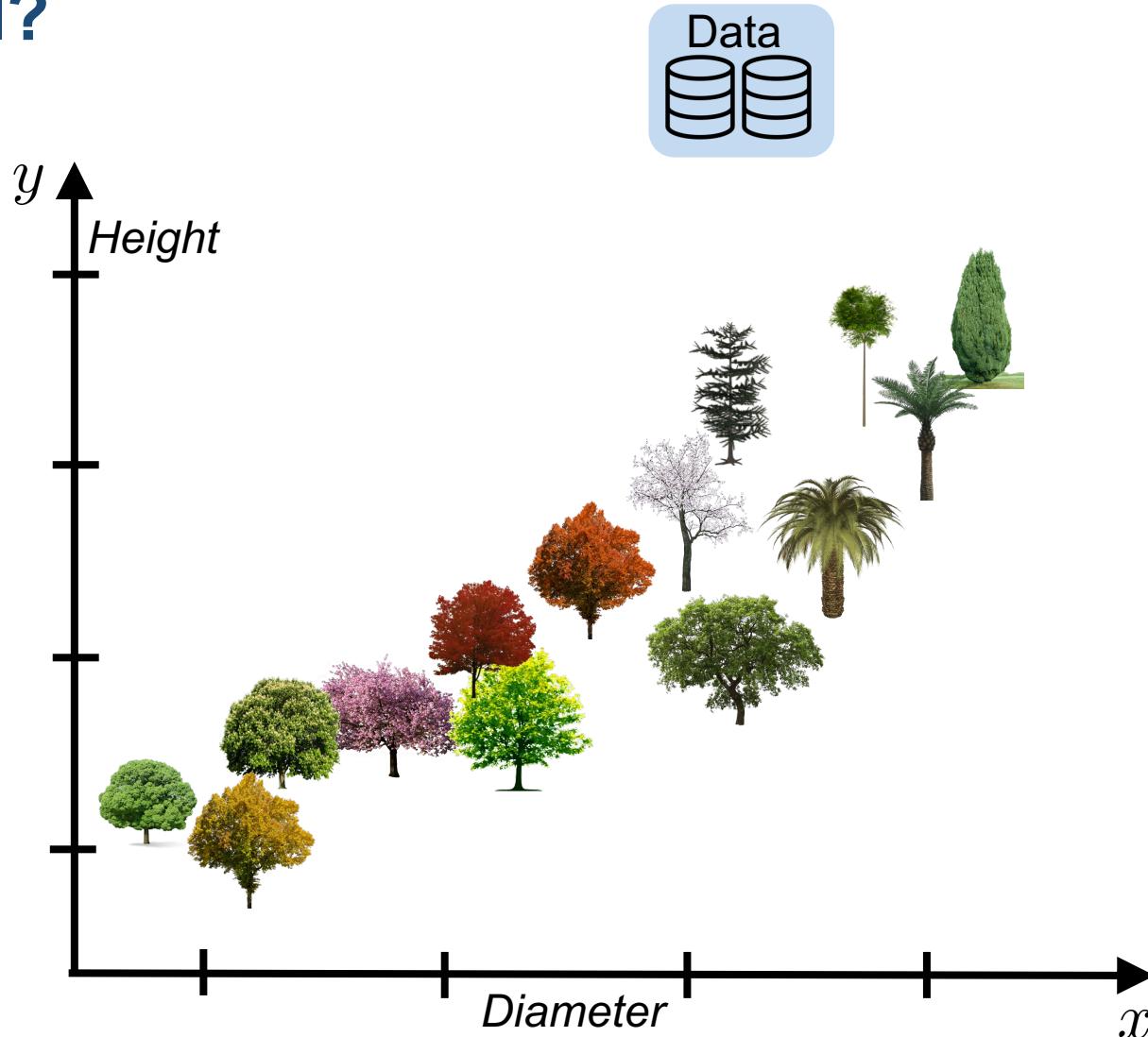
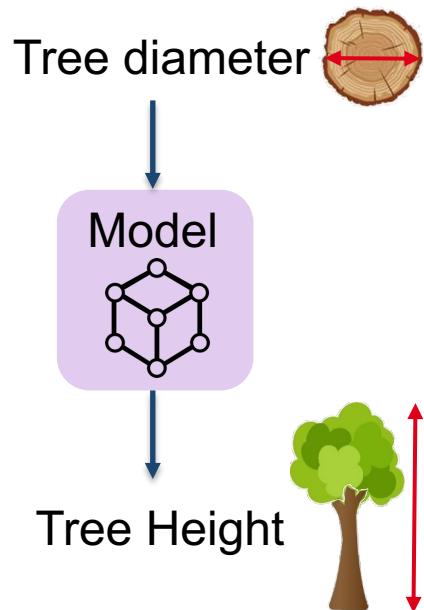
- Examples:
  - mean localization error in degrees*
  - mean classification error*
  - mean signal-to-error ratio in dB*
  - negative likelihood of observed data*
- The total loss  $L$  is designed based on the **task**, the **data**, and the chosen family of **models**. It measures the **fit** of  $m_{\hat{\theta}}$  for these data and task.
- Most modern machine learning **algorithms** can be interpreted as minimizing a loss, using **optimization**.



Animation by Andrew Ng

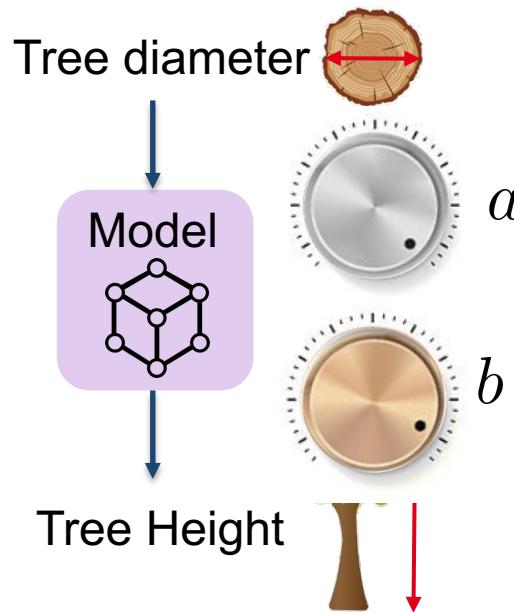
# What is a model?

*Example:*



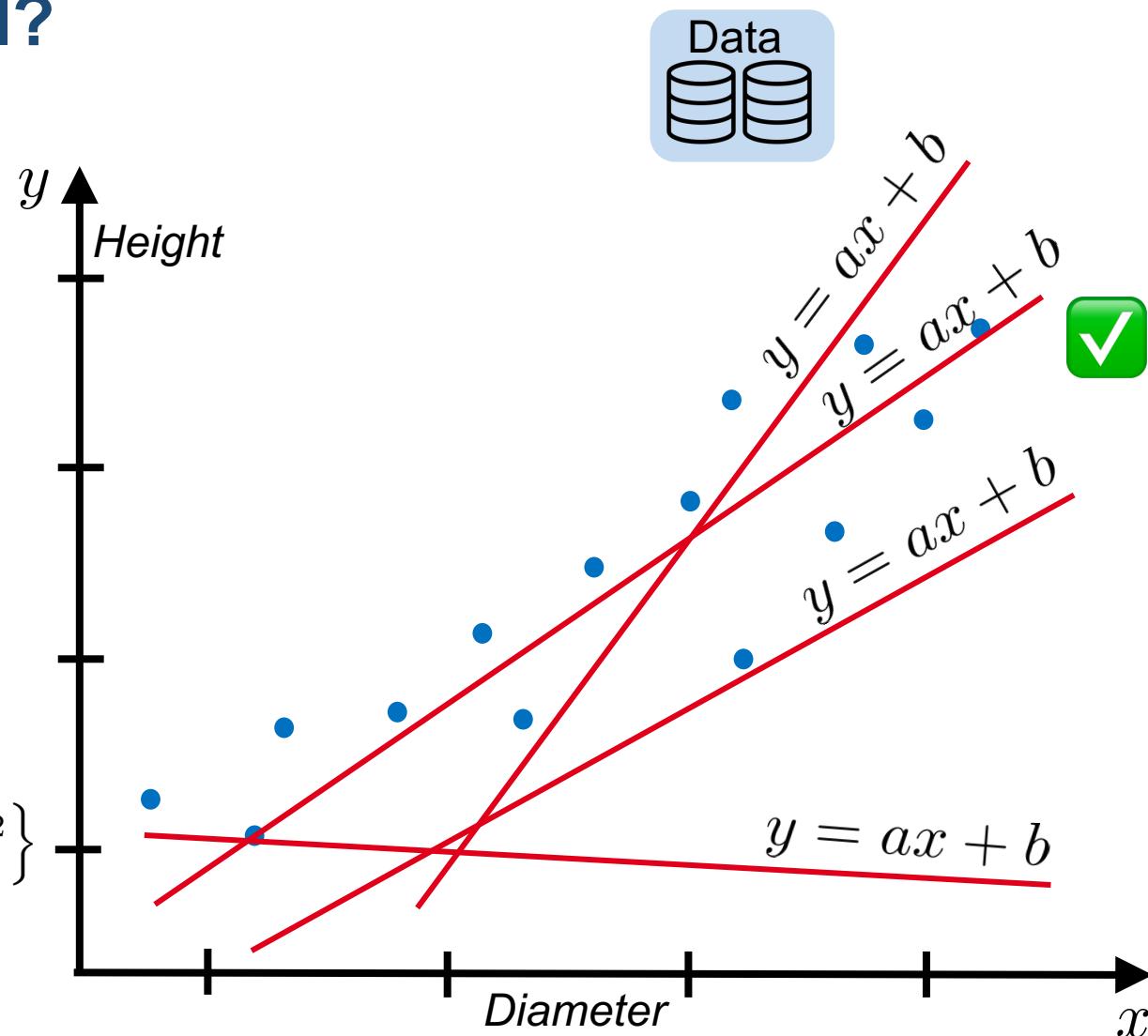
# What is a model?

**Example:**



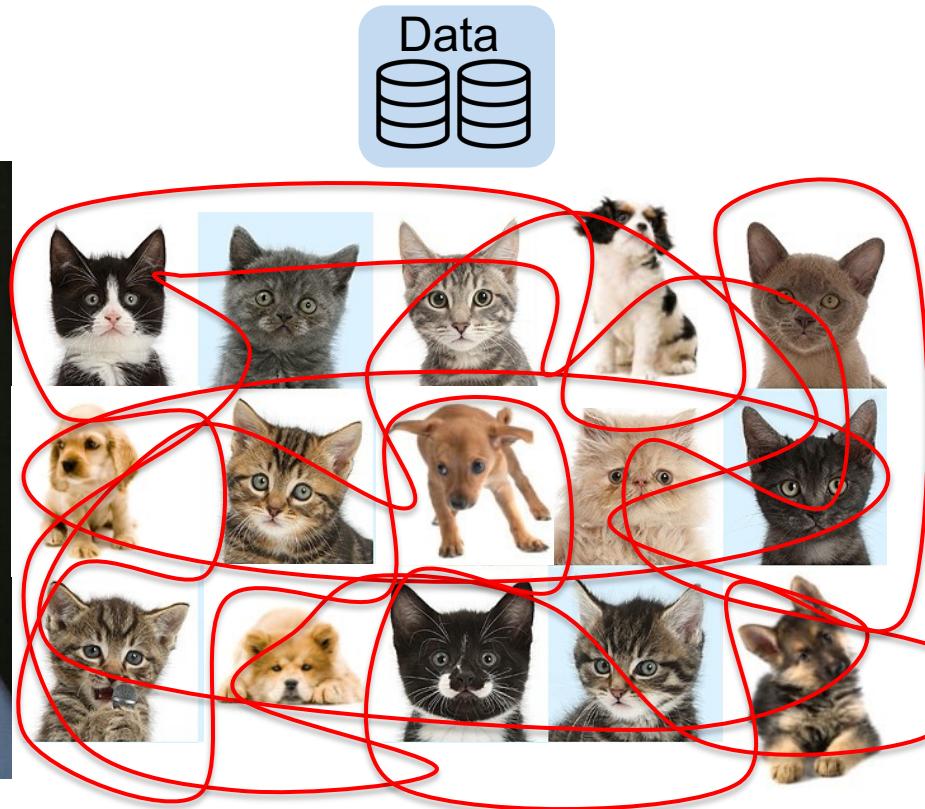
$$\underset{a,b}{\operatorname{argmin}} \left\{ \text{mean}_i (ax_i + b - y_i)^2 \right\}$$

« *The line is as close as possible to the points* »



# What is a model?

## Training stage

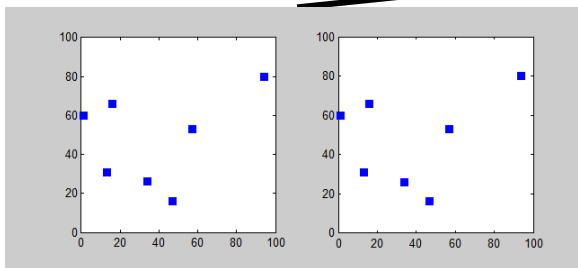


## Inference / evaluation stage

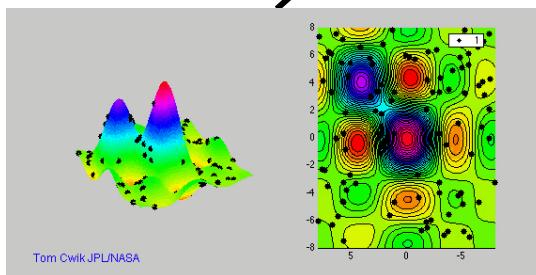


*It's a cat !*

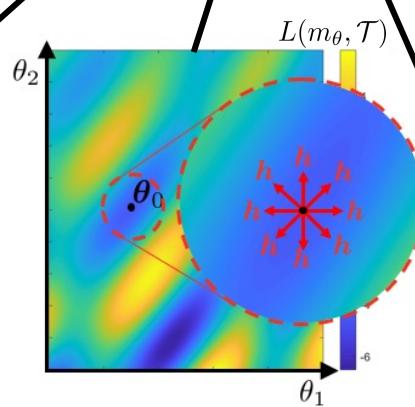
## Optimization Techniques



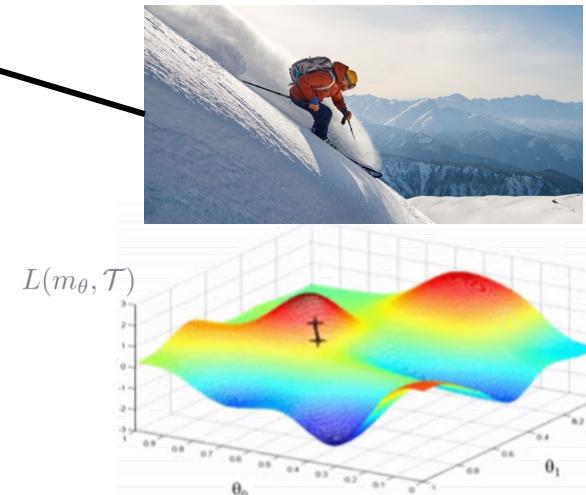
1. Brute Force



2. Population-Based



3. Zeroing gradient



5. Gradient Descent

$$\begin{aligned} \theta_1^{(i+1)} &= \operatorname{argmin}_{\theta_1} g(\theta_1, \theta_2^{(i)}, \dots, \theta_P^{(i)}) \\ \theta_2^{(i+1)} &= \operatorname{argmin}_{\theta_2} g(\theta_1^{(i+1)}, \theta_2, \dots, \theta_P^{(i)}) \\ &\vdots \\ \theta_P^{(i+1)} &= \operatorname{argmin}_{\theta_P} g(\theta_1^{(i+1)}, \theta_2^{(i+1)}, \dots, \theta_P) \end{aligned}$$

4. Alternate Minimization

## Artificial Intelligence

### “Applied” A.I.

#### Symbolic A.I.

#### Machine Learning

- **Algorithms** that (implicitly) build **models** from **data** to achieve **tasks**

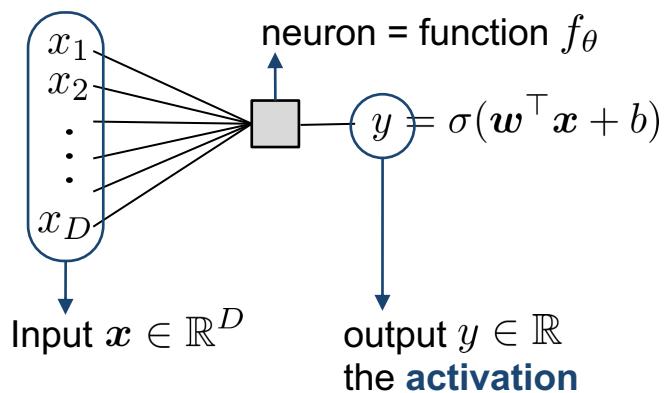
- Different **flavors** depending on available data and task (sup./unsup.)
- Main approach = **Model Fitting**: find a model that **minimizes a loss** within a **parameterized** family by **optimization**

#### Neural Networks

#### Deep Learning

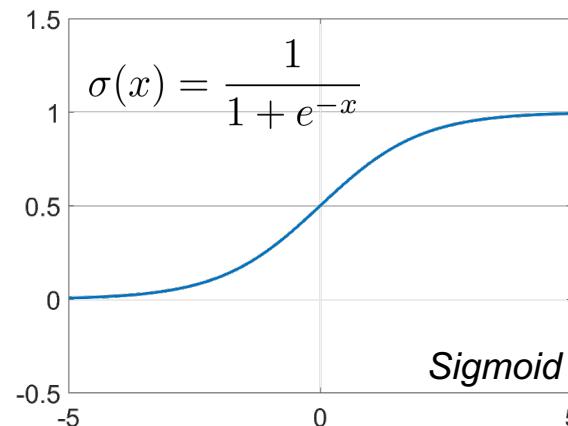
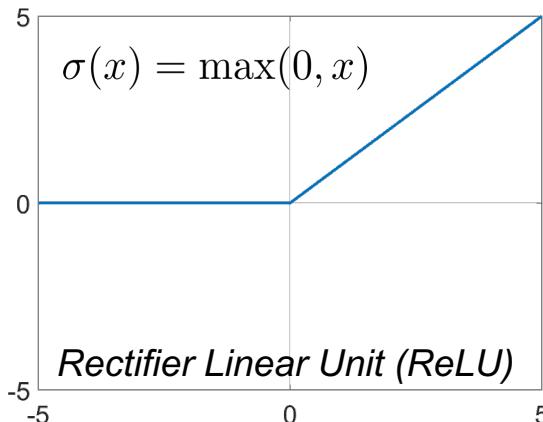
# Artificial Neural Networks

- **Artificial Neuron** = a multiple-input, single-output, **parametric, nonlinear function**



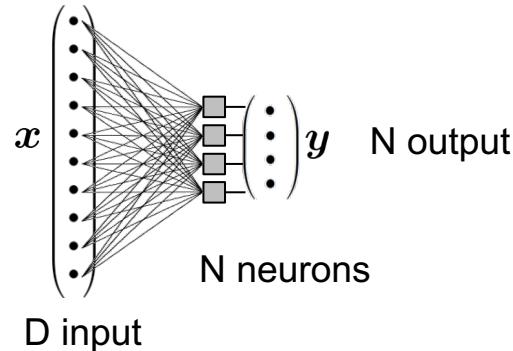
- $w^\top x = \langle w, x \rangle = \sum_{d=1}^D w_d x_d$ , the **dot product**
- $w = [w_1, \dots, w_D]^\top \in \mathbb{R}^D$ , the **weights**
- $b \in \mathbb{R}$ , the **bias**
- $w^\top x + b$ : the **pre-activation**
- $\theta = (w, b) \in \mathbb{R}^D \times \mathbb{R}$ , the **parameters**
- $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ , the **activation function** or **non-linearity**.

Examples of activation functions:



# Artificial Neural Networks

- **Simple perceptron:** multiple neurons attending the same input



- We have:  $y_n = \sigma(\mathbf{w}_n^\top \mathbf{x} + b_n)$
- Or in matrix form:  $\mathbf{y} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$ , where

$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_1^\top \\ \mathbf{w}_2^\top \\ \vdots \\ \mathbf{w}_N^\top \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{1,2} & \dots & w_{1,D} \\ w_{2,1} & w_{2,2} & \dots & w_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ w_{N,1} & w_{N,2} & \dots & w_{N,D} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix}, \quad \theta = (\mathbf{W}, \mathbf{b}) \in \mathbb{R}^{N \times D} \times \mathbb{R}^N$$

- Note: In this notation,  $\sigma$  is applied **elementwise** to a vector

## Artificial Neural Networks

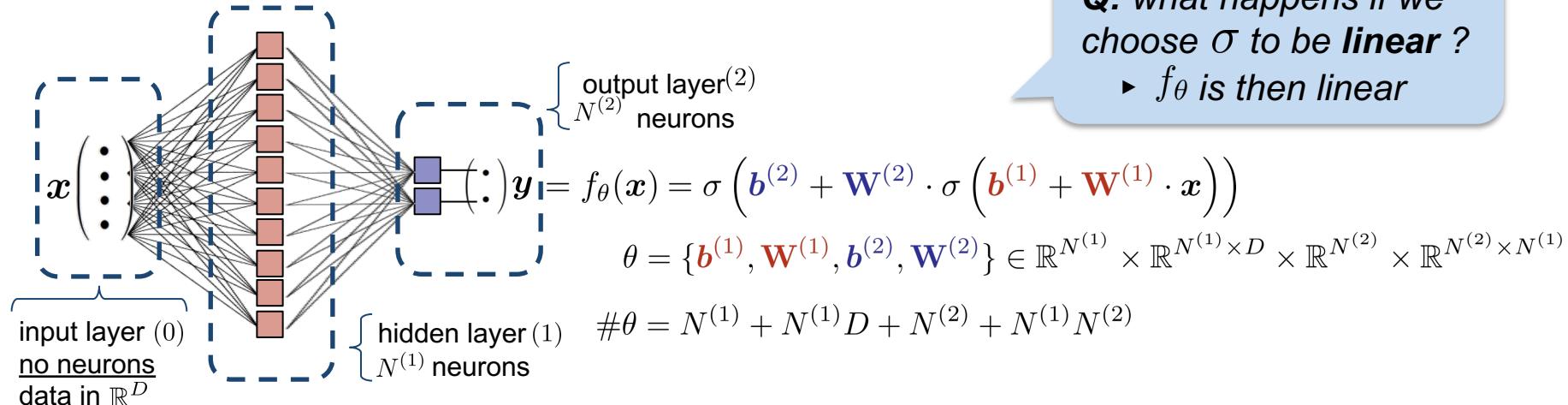
- Neurons can be **chained** together

$x$  →  $\boxed{(1)}$  →  $\boxed{(2)}$  →  $\boxed{(3)}$  →  $\boxed{(4)}$

$$y = f_{\theta}(x) = \sigma(b^{(4)} + w^{(4)} \cdot \sigma(b^{(3)} + w^{(3)} \cdot \sigma(b^{(2)} + w^{(2)} \cdot \sigma(b^{(1)} + w^{(1)} \cdot x))))$$

$$\theta = \{b^{(1)}, w^{(1)}, b^{(2)}, w^{(2)}, b^{(3)}, w^{(3)}, b^{(4)}, w^{(4)}\} \in \mathbb{R}^8$$

- **Simple Perceptron** (single hidden layer)



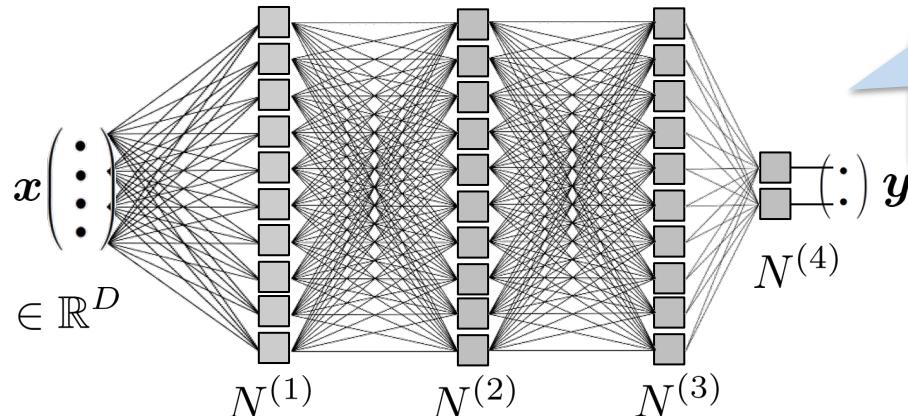
**Q:** what happens if we choose  $\sigma$  to be **linear** ?  
 ►  $f_{\theta}$  is then linear

**Universal Approximation Theorem** [Hornik, Stinchcomb, White, 1991]:  
 “A **single hidden layer** neural network with any “**sigmoid-like**” activation function and with a **linear output** unit can approximate any continuous function arbitrarily well, for **sufficiently large**  $N^{(1)}$ .”

... But the required  $N^{(1)}$  may be **exponential** in the input dimension  $D$ .

# Artificial Neural Networks

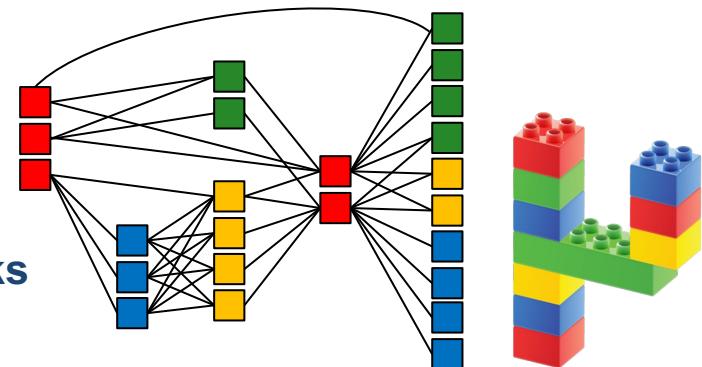
- Multilayer Perceptron (MLP) with 3 hidden layer (**Depth 4**)



Q:

Suppose  $D = N^{(1)} = N^{(2)} = N^{(3)} = N^{(4)} = N$ ,  
how many parameters?

$$\#\theta = 4(N + N^2) = \mathcal{O}(N^2)$$



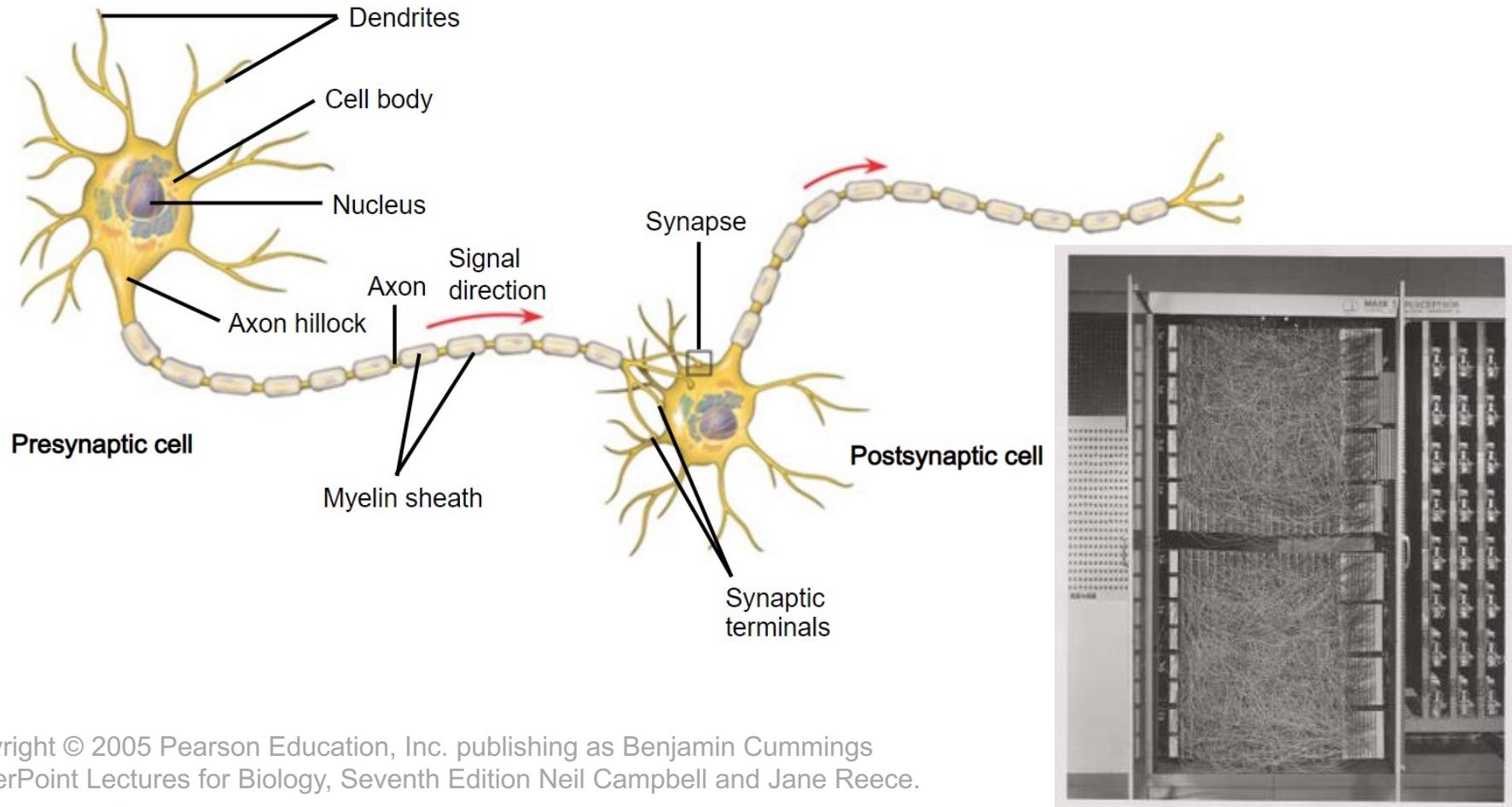
- Artificial neurons, as **elementary computing units**, can be combined in many different ways
- No cycle in the graph = **Feedforward Neural Networks**
- We call a given network of neurons an **architecture**
- Neural Networks form a very flexible **class of parameterized families** of **non-linear functions**:  $\{\mathcal{F}_i\}_i$  where  $\mathcal{F}_i = \{f_\theta^i\}_\theta$

Artificial neural networks with 2 or more hidden layers are called **Deep Neural Networks (DNNs)**

## Why Deep?

### 1) Biological inspiration

Frank Rosenblatt. « The perceptron: A probabilistic model for information storage and organization in the brain ». *Psychological Review* (1958).

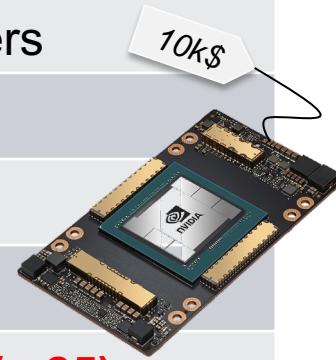


Copyright © 2005 Pearson Education, Inc. publishing as Benjamin Cummings  
PowerPoint Lectures for Biology, Seventh Edition Neil Campbell and Jane Reece.

## Why Deep?

### 1) Biological inspiration

Human Brain	Artificial Neural Network
~ 86 billion neurons	100k - 1 billion neurons
~ 7,000 synapse connections per neuron (~600 trillion connections)	3 - 1,000 connections per neuron 1 million - 1 trillion parameters
Massively parallel	(Mostly) sequential
Asynchronous	Synchronous
Very plastic architecture	(Mostly) fixed architecture
12 W of power	Nvidia A100 GPU : 300 W* ( <b>x 25</b> )
Biological neuron = extremely complex** and not fully understood	Artificial Neuron = a weighted sum and a threshold.
3.7B year evolution + experiencing life	Designed by humans, trained on Internet



10k\$

\* - A chat session with chatGPT mobilizes roughly a full A100. Training GPT-3 is estimated to have taken ~95 years of A100.

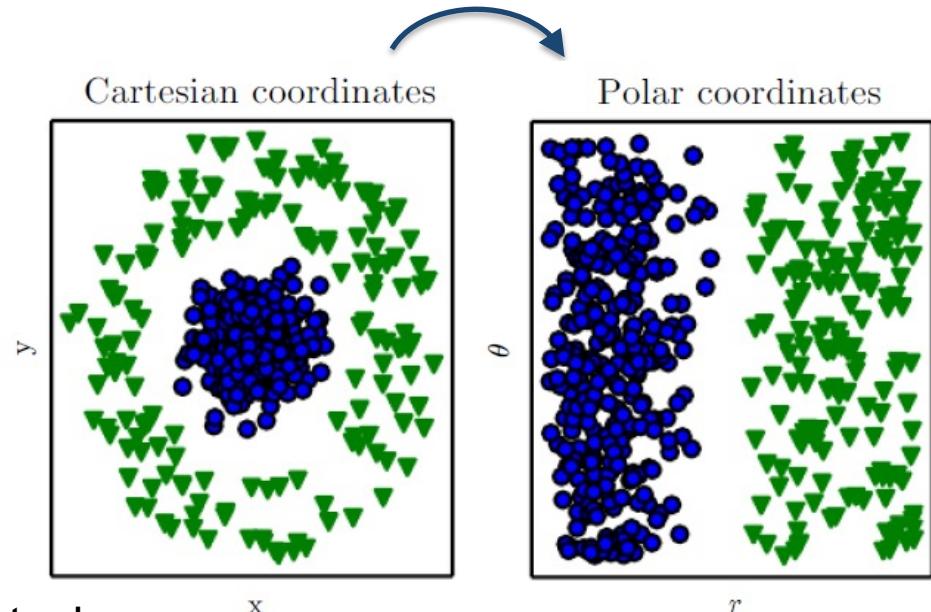
- GPT-3 has ~135 billion parameters and roughly ~0.25 billion “neurons”.

\*\* [Beniagiev et al. 2021]: a single bio-neuron is well approximated by a 1000-neuron ANN of depth 5-8

## Why Deep?

### 2) Bypassing feature engineering

- “Classical” machine learning relies heavily on **hand-crafted features**
- Such pre-processing allows successfully applying (combinations of) **simple / low-parametric** models
- But **difficult to do by hand** for many tasks

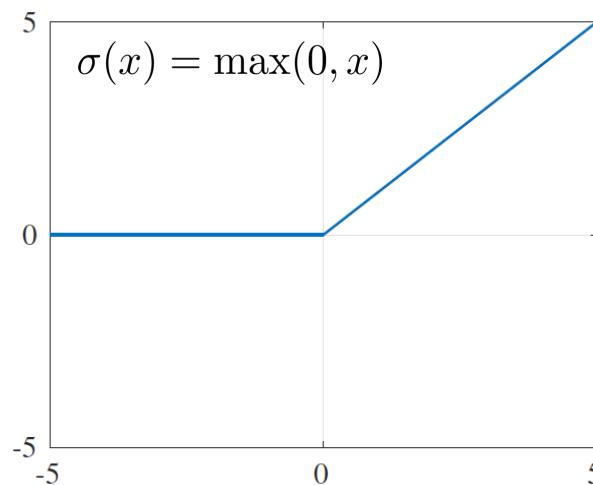


- **Sufficiently large DNNs** (and in particular deep **convolutional nets**) tend to **learn features** automatically. → The **deeper** the **higher** level.
- Can be successfully applied to **raw signals** (e.g. image pixels, audio samples)

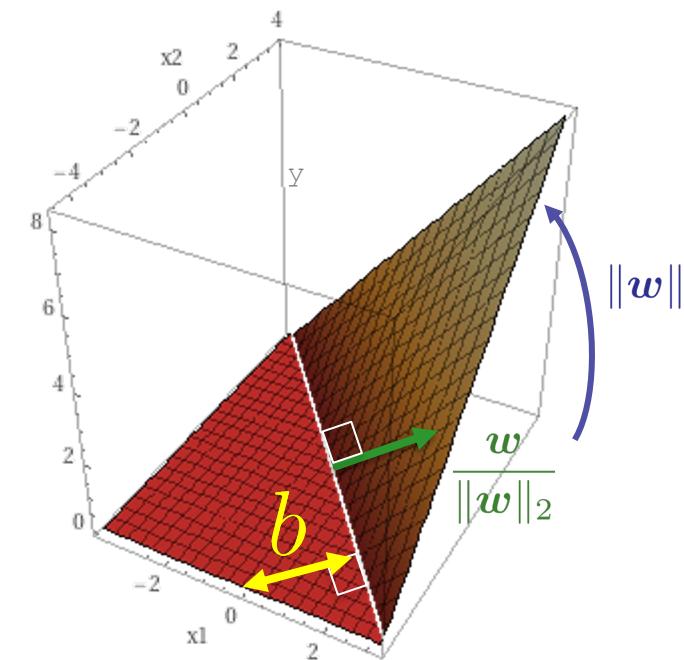
## Why Deep?

### 3) The “Origami Effect”

- ReLU activation:



- For a 2D input:  
 $y = \sigma(\mathbf{w}^\top \mathbf{x} + b)$

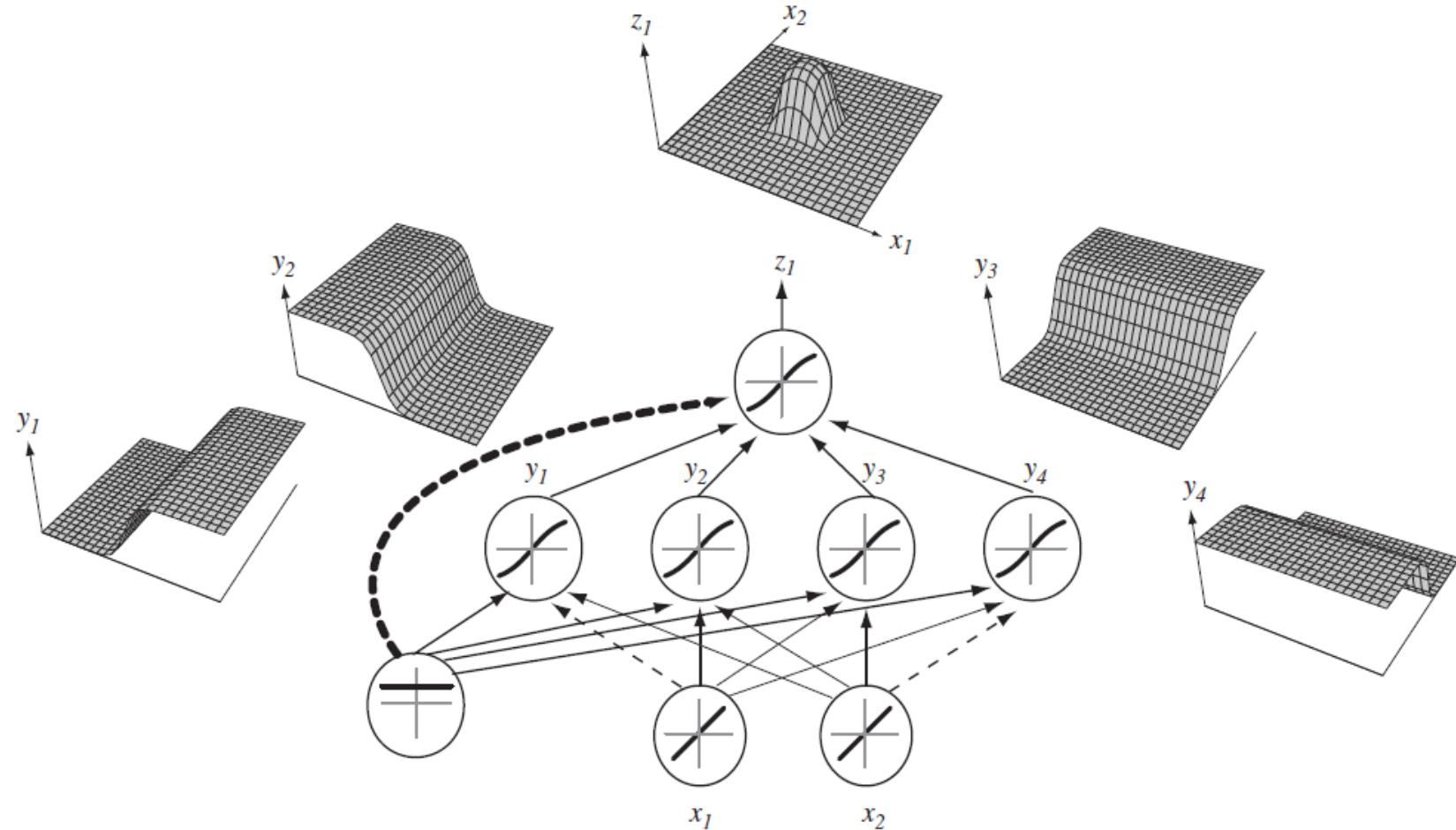


**Theorem** [G. Montufar et al., 2014]: *the max. number of linear regions modeled by a piecewise linear network (i.e., a network with ReLU neurons) with  $D$  inputs,  $L$  layers, and  $N$  units per layer is in the order of*

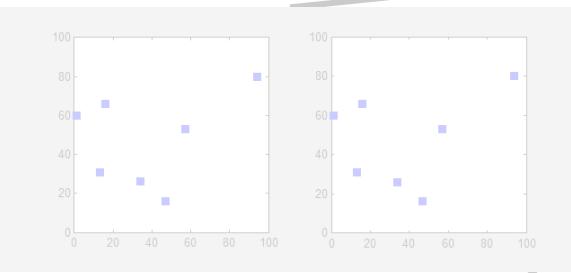
$$N^D \left(\frac{N}{D}\right)^{D(L-2)}, \text{ i.e., the model capacity is **exponential in the depth  $L$** ,}\\ \text{i.e., in the **model size** (recall it is } O(LN^2) \text{ )}$$

# Why Deep?

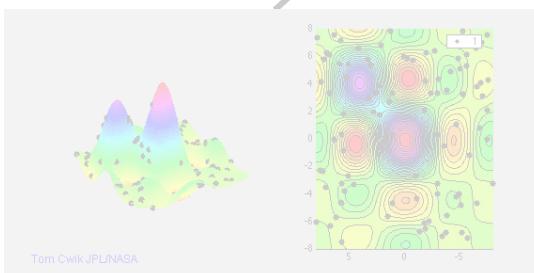
## 3) The “Origami Effect”: Illustration with Sigmoids



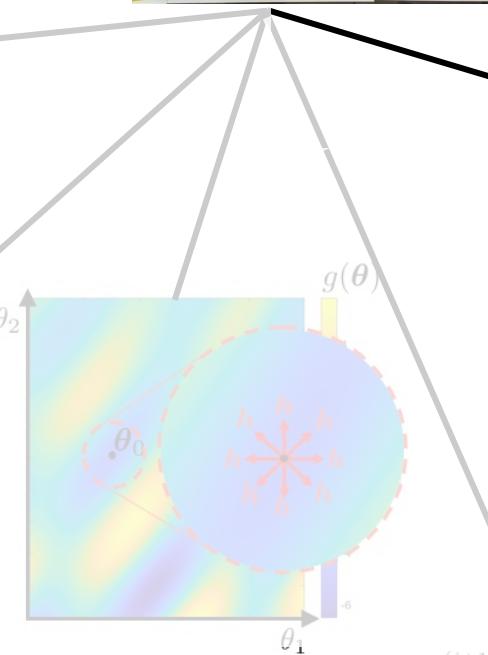
## How to train a DNN?



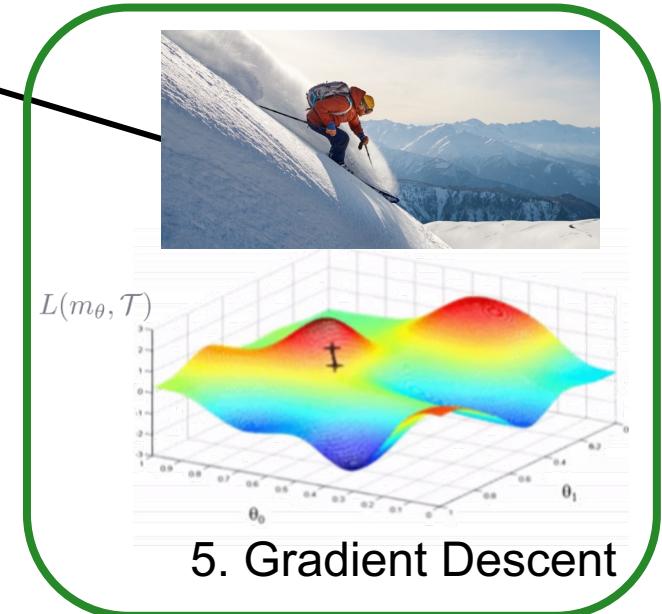
1. Brute Force



2. Population-Based



3. Zeroing gradient



**This one !**

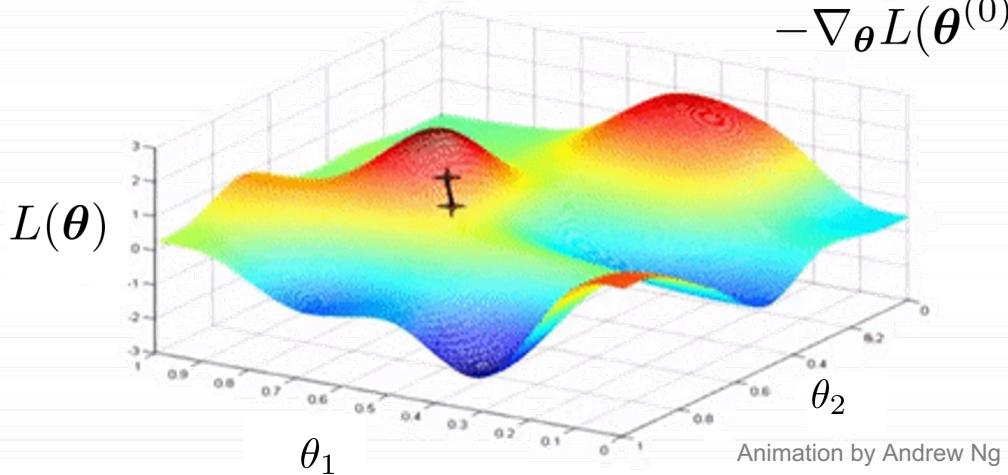
$$\begin{aligned} \theta_1^{(i+1)} &= \operatorname{argmin}_{\theta_1} g(\theta_1, \theta_2^{(i)}, \dots, \theta_P^{(i)}) \\ \theta_2^{(i+1)} &= \operatorname{argmin}_{\theta_2} g(\theta_1^{(i+1)}, \theta_2, \dots, \theta_P^{(i)}) \\ &\vdots \\ \theta_P^{(i+1)} &= \operatorname{argmin}_{\theta_P} g(\theta_1^{(i+1)}, \theta_2^{(i+1)}, \dots, \theta_P) \end{aligned}$$

4. Alternate Minimization

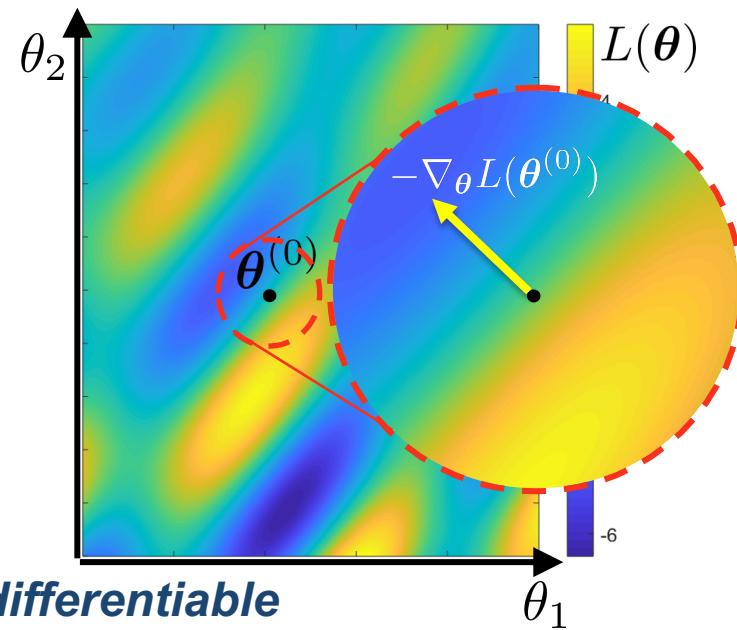
# Gradient Descent

## Intuition:

- Start from an initial **parameter vector**  $\theta^{(0)} \in \mathbb{R}^P$
- From here, follow the **direction of steepest descent**
- Stop when things look **flat**



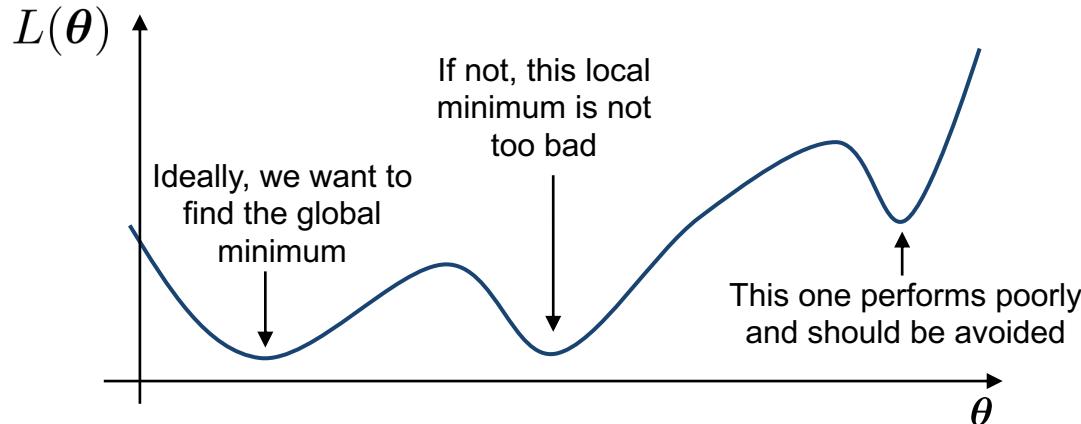
$$-\nabla_{\theta} L(\theta^{(0)}) = - \left[ \frac{\partial L}{\partial \theta_1} \Big|_{\theta_1^{(0)}}, \dots, \frac{\partial L}{\partial \theta_D} \Big|_{\theta_D^{(0)}} \right]^T$$



- The updates are: 
$$\boxed{\theta^{(i+1)} \leftarrow \theta^{(i)} - \epsilon \nabla_{\theta} L(\theta^{(i)})}$$
- Requires the function to be (almost everywhere) **differentiable**

# Gradient Descent

→ **The algorithm converges to local minima under mild assumptions** 😊



- **Spurious local minima** cannot always be avoided
- Many variants have been derived to limit them, and to speed up convergence

- The **gradient-step  $\epsilon$** , also called **learning rate**, is a critical **hyper-parameter**

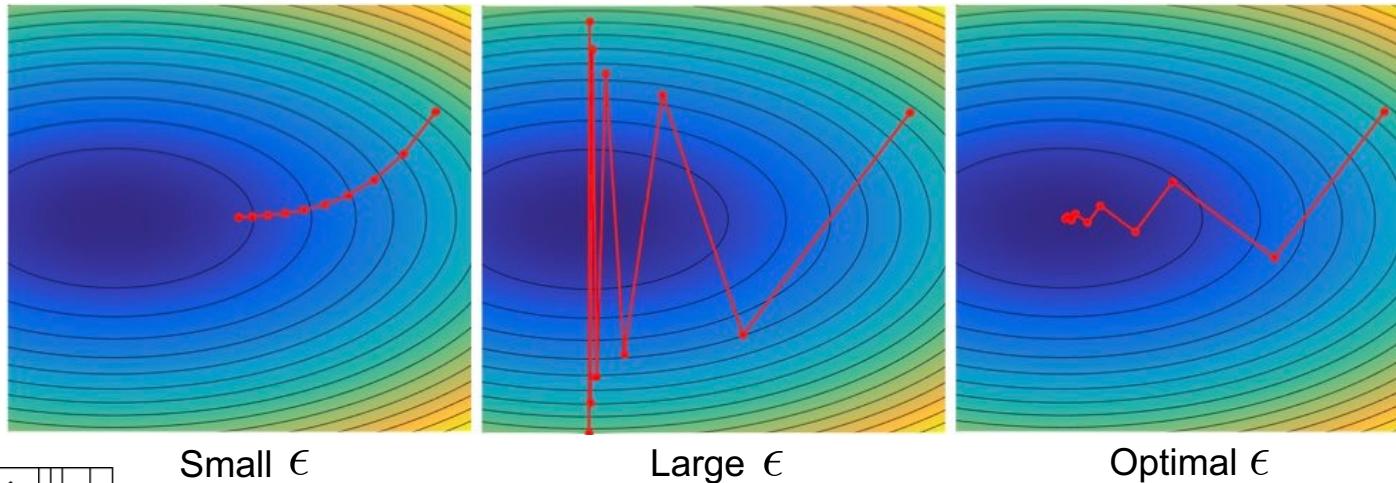
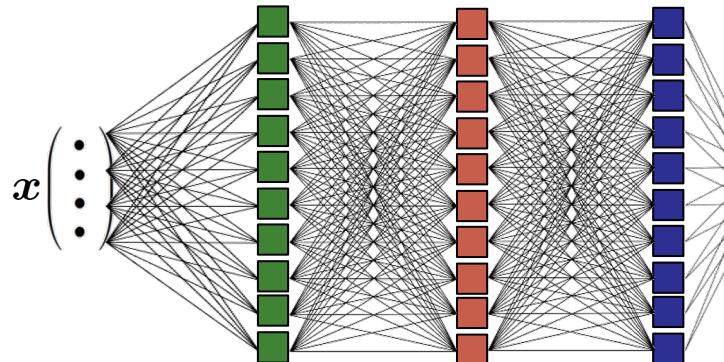


Image by  
Gabriel Peyre

**Problem:** How to efficiently compute the gradient of a DNN with respect to its parameters?



Compute  $\nabla_{\theta} L(\theta, \mathcal{T})$  ?!



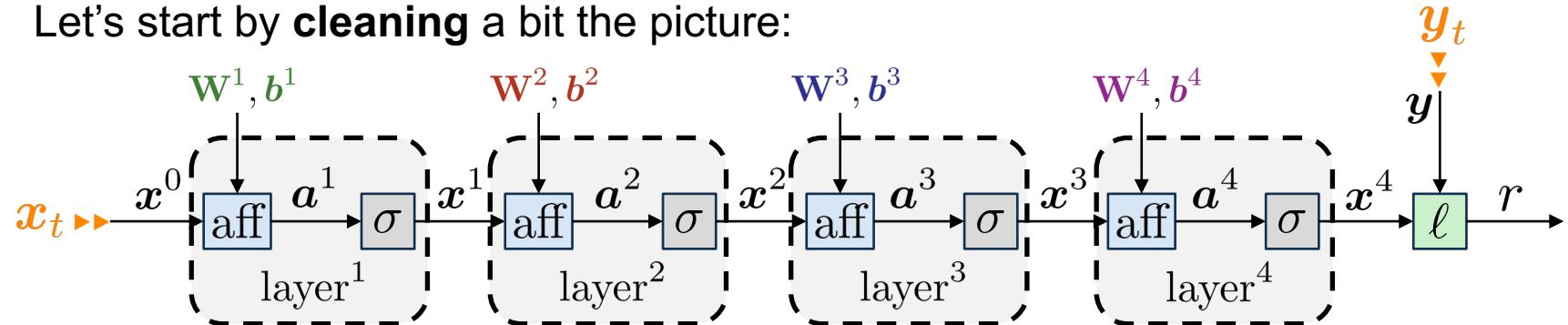
- Parameters:  $\theta = [\theta^i]_{i=1}^4 = [(b^i, W^i)]_{i=1}^4$
- Training set:  $\mathcal{T} = \{(x_t, y_t)\}_{t=1}^T$
- Total Loss function:  $L(\theta, \mathcal{T}) = \sum_{t=1}^T \ell(\text{dnn}_{\theta}(x_t), y_t)$

## The *Backpropagation* Algorithm

- Uses the **chain rule** to compute the gradient layer by layer
- **Forward step:** send a data point ( $x_t, y_t$ ) through the network's layers
- **Backward step:** compute the loss gradient at this point by recursively multiplying intermediate gradient values

# The Backpropagation Algorithm

Let's start by **cleaning** a bit the picture:

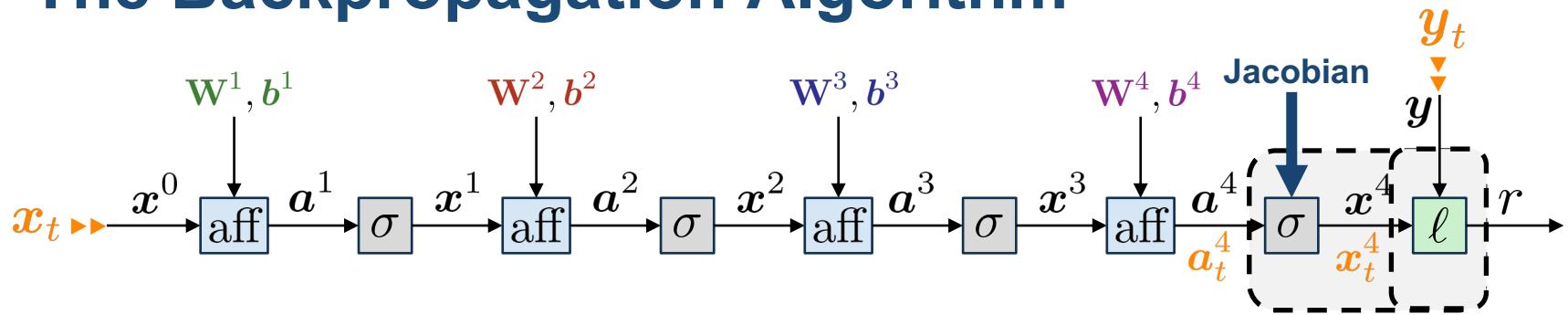


- $\mathbf{a}^i = \text{aff}_{\theta^i}(\mathbf{x}^{i-1}) = \mathbf{W}^i \mathbf{x}^{i-1} + \mathbf{b}^i$  are the **pre-activations**
- $\mathbf{x}^i = \sigma(\mathbf{a}^i) = \sigma(\text{aff}_{\theta^i}(\mathbf{x}^{i-1})) = \text{layer}^i(\mathbf{x}^{i-1})$  are the **activations**
- The **loss**  $\ell$  can be viewed as another layer, with **real output**  $r$  (the “*residual*”)
- Again, we have:  $\nabla_{\theta} L(\theta, \mathcal{T}) = \frac{1}{T} \sum_{t=1}^T \nabla_{\theta} \ell(\text{dnn}_{\theta}(\mathbf{x}_t), \mathbf{y}_t)$
- It's enough to calculate the gradient of the loss at **one sample**  $(\mathbf{x}_t, \mathbf{y}_t)$ :

$$\mathbf{G}_{\theta} \stackrel{\text{def}}{=} \nabla_{\theta} \ell(\text{dnn}_{\theta}(\mathbf{x}_t), \mathbf{y}_t)$$

and then sum.

# The Backpropagation Algorithm



0) We start by  $G_{x^4} = \frac{\partial r}{\partial x^4} \Big|_{x_t^4}^\top = \nabla_{x^4} \ell(x_t^4, y_t)$ .

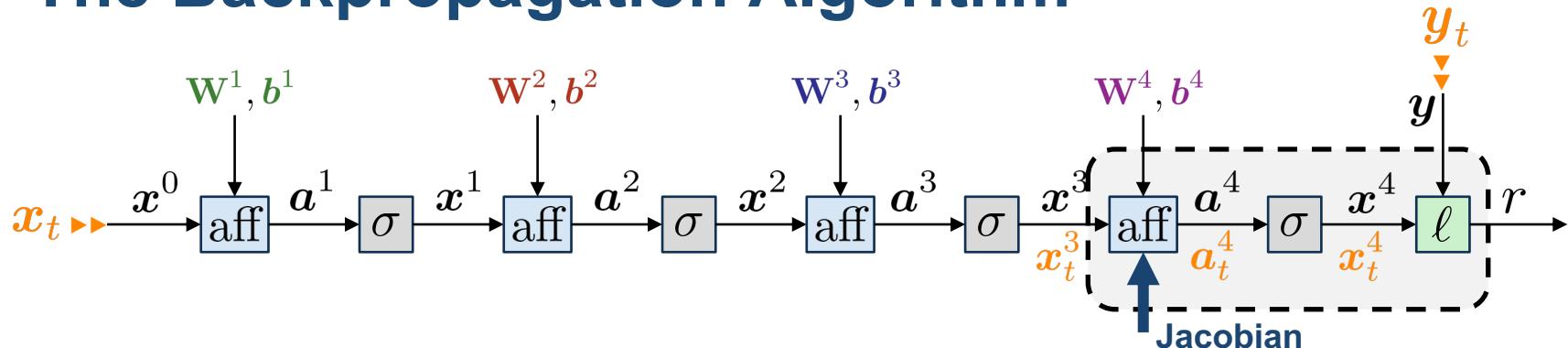
For example, for the L2 loss  $\ell(x^4, y_t) = \|x^4 - y_t\|_2^2$ , we have  $G_{x^4} = 2(x_t^4 - y_t)$ , the **difference** between the network **prediction** and the **target**.

1) Then,  $G_{a^4} = \frac{\partial r}{\partial a^4} \Big|_{a_t^4}^\top = \left( \frac{\partial r}{\partial x^4} \Big|_{x_t^4} \times \frac{\partial x^4}{\partial a^4} \Big|_{a_t^4} \right)^\top = \boxed{\frac{\partial x^4}{\partial a^4} \Big|_{a_t^4}^\top} \times G_{x^4}$

$$= \sigma'(a_t^4) \odot G_{x^4}.$$

diag[ $\sigma'(a_t^4)$ ] !

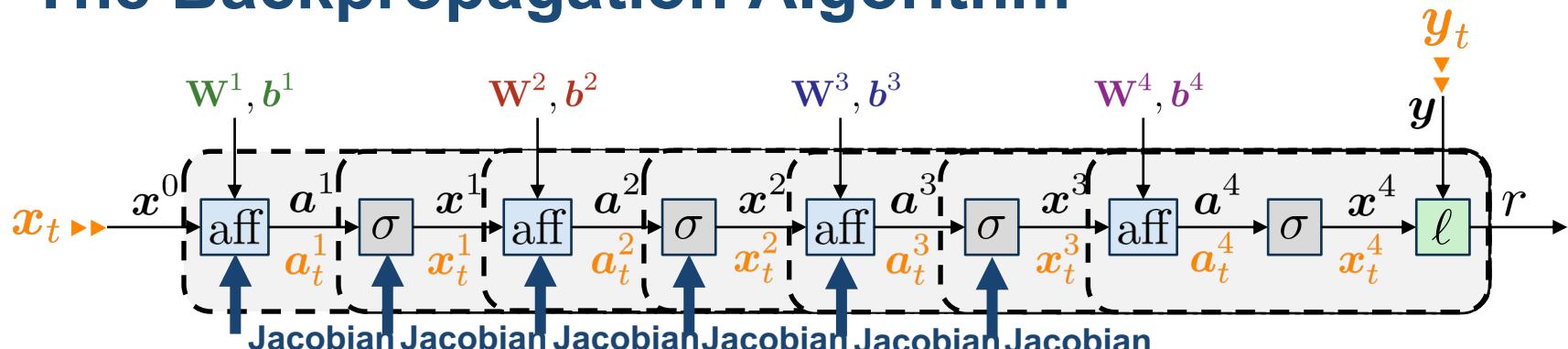
# The Backpropagation Algorithm



2) Then:

- $G_{x^3} = \frac{\partial r}{\partial x^3} \Big|_{x_t^3}^\top = \left( \frac{\partial r}{\partial a^4} \Big|_{a_t^4} \times \frac{\partial a^4}{\partial x^3} \Big|_{x_t^3} \right)^\top = \boxed{\frac{\partial a^4}{\partial x^3} \Big|_{x_t^3}^\top} \times G_{a^4} = \boxed{W^4^\top G_{a^4}}.$
- $G_{b^4} = \frac{\partial r}{\partial b^4} \Big|_{b^4}^\top = \left( \frac{\partial r}{\partial a^4} \Big|_{a_t^4} \times \frac{\partial a^4}{\partial b^4} \Big|_{b^4} \right)^\top = \boxed{\frac{\partial a^4}{\partial b^4} \Big|_{b^4}^\top} \times G_{a^4} = \boxed{G_{a^4}}.$
- $G_{w_d^4} = \frac{\partial r}{\partial w_d^4} \Big|_{w_d^4}^\top = \left( \frac{\partial r}{\partial a^4} \Big|_{a_t^4} \times \frac{\partial a^4}{\partial w_d^4} \Big|_{w_d^4} \right)^\top = \boxed{\frac{\partial a^4}{\partial w_d^4} \Big|_{w_d^4}^\top} \times G_{a^4} = \boxed{x_{t,d}^3 G_{a^4}}.$

# The Backpropagation Algorithm



And so on...

$$3) \quad G_{a^3} = \sigma'(\mathbf{a}_t^3) \odot G_{x^3} . \quad 5) \quad G_{a^2} = \sigma'(\mathbf{a}_t^2) \odot G_{x^2} . \quad 7) \quad G_{a^1} = \sigma'(\mathbf{a}_t^1) \odot G_{x^1} .$$

$$4) \quad G_{x^2} = \mathbf{W}^{3\top} G_{a^3} . \quad 6) \quad G_{x^1} = \mathbf{W}^{2\top} G_{a^2} . \quad 8) \quad G_{x^0} = \mathbf{W}^{1\top} G_{a^1} .$$

$$G_{b^3} = G_{a^3} .$$

$$G_{b^2} = G_{a^2} .$$

$$G_{b^1} = G_{a^1} .$$

$$G_{w_d^3} = x_{t,d}^2 G_{a^3} .$$

$$G_{w_d^2} = x_{t,d}^1 G_{a^2} .$$

$$G_{w_d^1} = x_{t,d}^0 G_{a^1} .$$

Putting everything together, we have for example:

$$G_{b^1} = \sigma'(\mathbf{a}_t^1) \odot \mathbf{W}^{2\top} \sigma'(\mathbf{a}_t^2) \odot \mathbf{W}^{3\top} \sigma'(\mathbf{a}_t^3) \odot \mathbf{W}^{4\top} \sigma'(\mathbf{a}_t^4) \odot \nabla_{\mathbf{x}^4} \ell(\mathbf{x}_t^4, \mathbf{y}_t) .$$

# Stochastic Gradient Descent (SGD) and Batching

- We train DNNs using so called ***Empirical Risk Minimization***:

$$L(\boldsymbol{\theta}, \mathcal{T}) = \frac{1}{T} \sum_{t=1}^T \ell(\text{dnn}_{\boldsymbol{\theta}}(\mathbf{x}_t), \mathbf{y}_t) \approx \mathbb{E}_{\mathbf{X}, \mathbf{Y}} \{ \ell(\text{dnn}_{\boldsymbol{\theta}}(\mathbf{X}), \mathbf{Y}) \}$$

- We compute the gradient of the **total loss** by summing the gradients of the **loss** at individual data samples:

$$\nabla_{\boldsymbol{\theta}} L(\text{dnn}_{\boldsymbol{\theta}}, \mathcal{T}) = \frac{1}{T} \sum_{t=1}^T \nabla_{\boldsymbol{\theta}} \ell(\text{dnn}_{\boldsymbol{\theta}}(\mathbf{x}_t), \mathbf{y}_t)$$

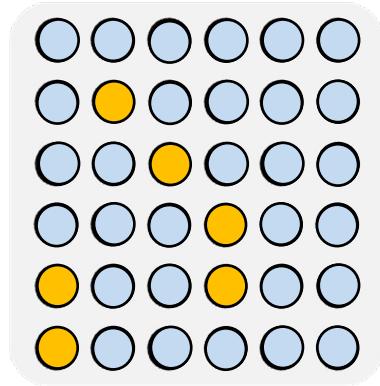
- But doing so across the **entire dataset** (e.g.: 1 million images) for **every gradient step** would be very expansive.

# Stochastic Gradient Descent (SGD) and Batching

- At each iteration ( $i$ ), compute the gradient over a **random subset**  $\mathcal{T}^{(i)} \subseteq \mathcal{T}$  and perform one step of gradient descent:

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \epsilon \cdot \frac{1}{|\mathcal{T}^{(i)}|} \sum_{(\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{T}^{(i)}} \nabla_{\theta} \ell(\text{dnn}_{\theta}(\mathbf{x}_t), \mathbf{y}_t)$$

- At next iteration, pick another (disjoint) random subset
- When the entire dataset has passed, start over again
- Each  $\mathcal{T}^{(i)}$  is called a **minibatch**
- Each pass over the entire dataset is called an **epoch**



Splitting the training set into  $B$  minibatches:

- Reduces the computation cost of one gradient by a factor of  $B$
- Increases the **standard deviation** on the gradient estimate by a factor of  $\sqrt{B}$  only.

*More iterations but fewer epochs  
= less total computation*

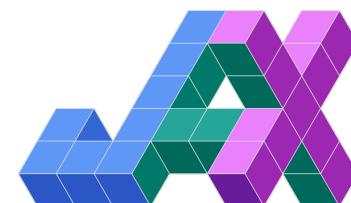
*+ compute gradients of 1 batch in parallel via GPU*

# The PyTorch framework



**GOOD NEWS:** You (probably) won't ever need to implement backpropagation or SGD yourself! 😊

- PyTorch: opensource Python library designed to easily **design, train** and **test** neural networks
- Initially developed by Facebook (Meta), based on Torch. **Constantly evolving** thanks to a broad community
- Seamlessly manipulate all the objects we have seen thanks to differential programming : **data, batches, parameters, layers, optimizers, loss...**
- Includes support for **GPU** and a C++ interface.
- Competing frameworks: **TensorFlow, Jax**



# Outline

---

## 1. Introduction ✓

AI? Machine Learning? Neural Networks?

## 2. Machine Learning ✓

Definition, Characteristics, Types, Algorithms

## 3. Deep Neural Networks ✓

Definition, Motivation, How to Train Them

## 4. Limits of Machine Learning

Overfitting, Data Biases

# Outline

---

## 1. Introduction

AI? Machine Learning? Neural Networks?

## 2. Machine Learning

Definition, Characteristics, Types, Algorithms

## 3. Deep Neural Networks

Definition, Motivation, How to Train Them

## 4. Limits of Machine Learning

Overfitting, Data Biases

## Scenario:

- Imagine we have a training dataset  $\mathcal{T}$  containing 10,000 images with labels (supervised learning)
- We train a machine learning model to perform the multi-class classification task
- We get nearly **perfect results** on these 10,000 images, e.g., 99.9% of correct classification
- However, when we run the model on new images, the results are **awful**, i.e., close to random.

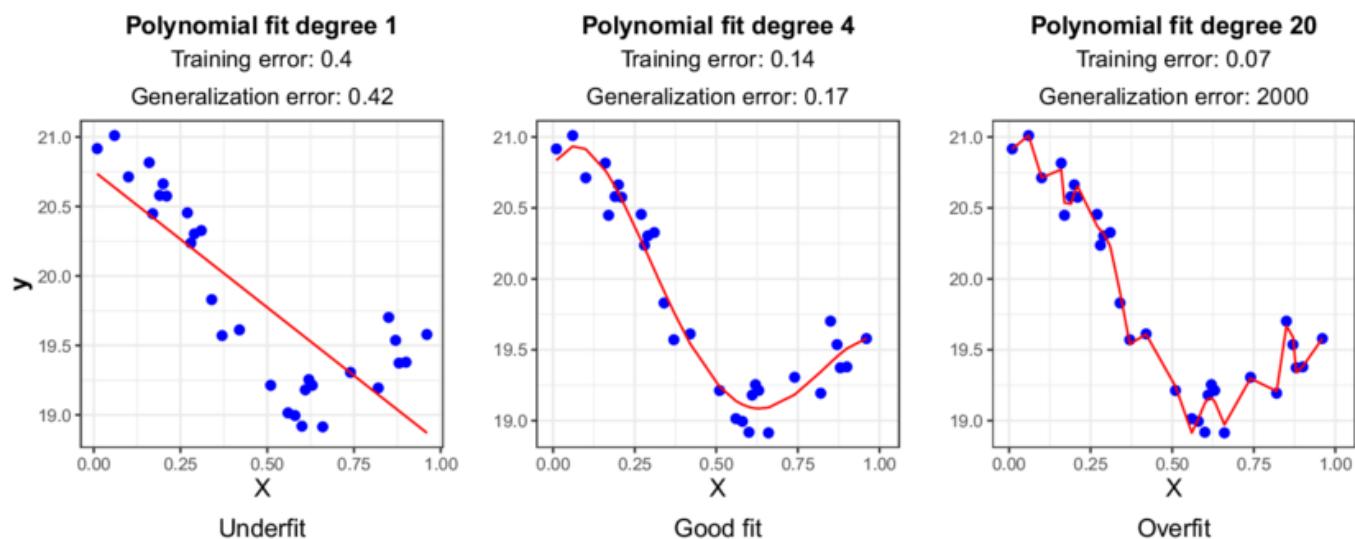
*What's going on?*

# Overfitting

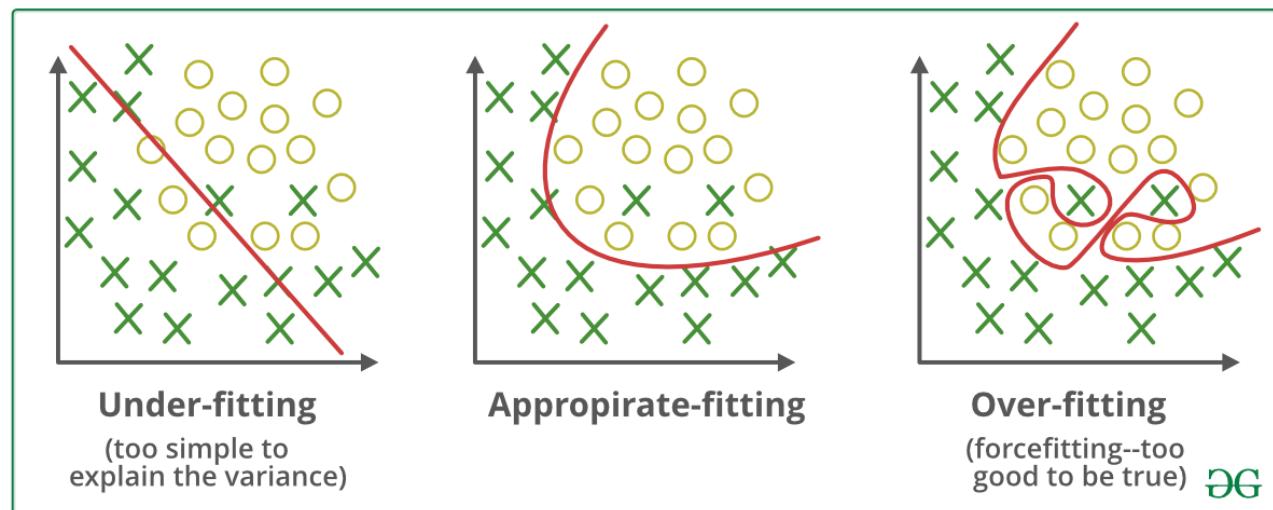
- Our algorithm is guilty of **overfitting**
- Instead of learning general features to classify the images, it **learned by heart** all the images in our training dataset!
- We sometimes have **millions** or **billions** of parameters in a model => it has the **capacity** to **store/encode** large amounts of data
- This may even happen for models of relatively small capacity, if the **amount of training data** is insufficient.

## Overfitting

- Ex: polynomial regression



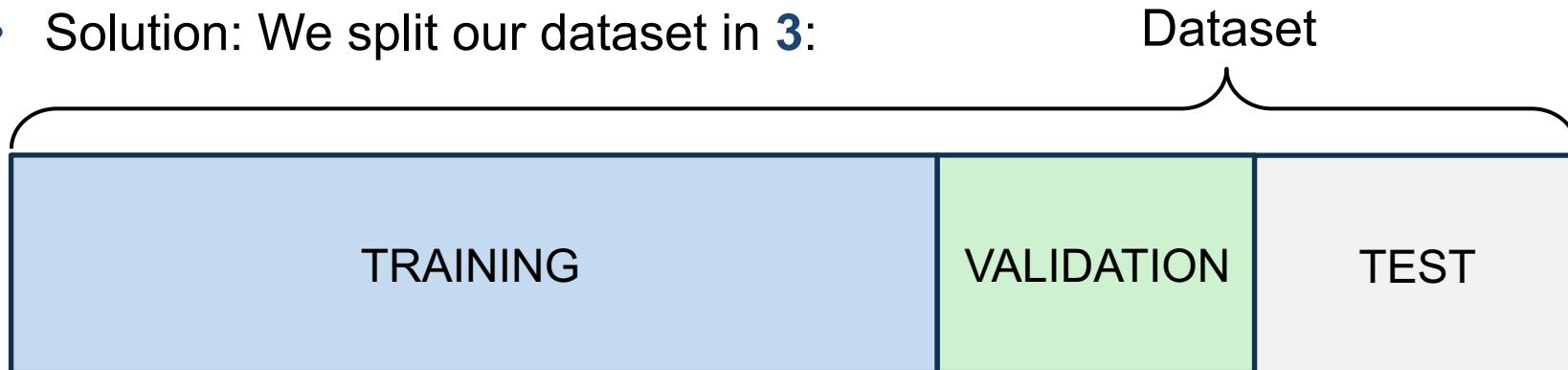
- Ex: binary classification



DG

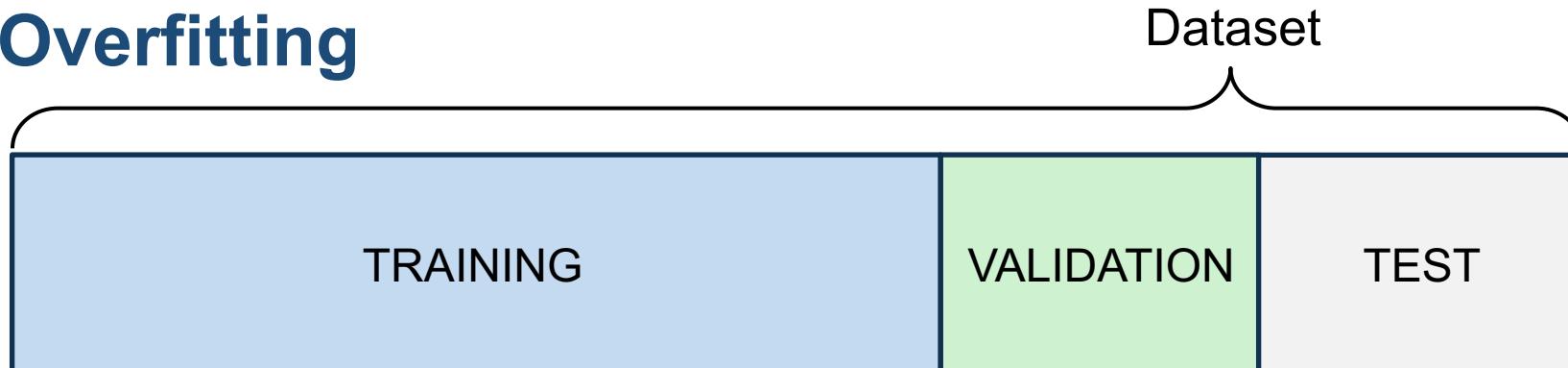
## Overfitting

- In supervised learning, we are not really interested in a model that works perfectly on our training set. **We already have the answers anyway**, by definition of a (supervised) training set !
- We want a model that **generalizes** to **unseen** data
- Solution: We split our dataset in **3**:



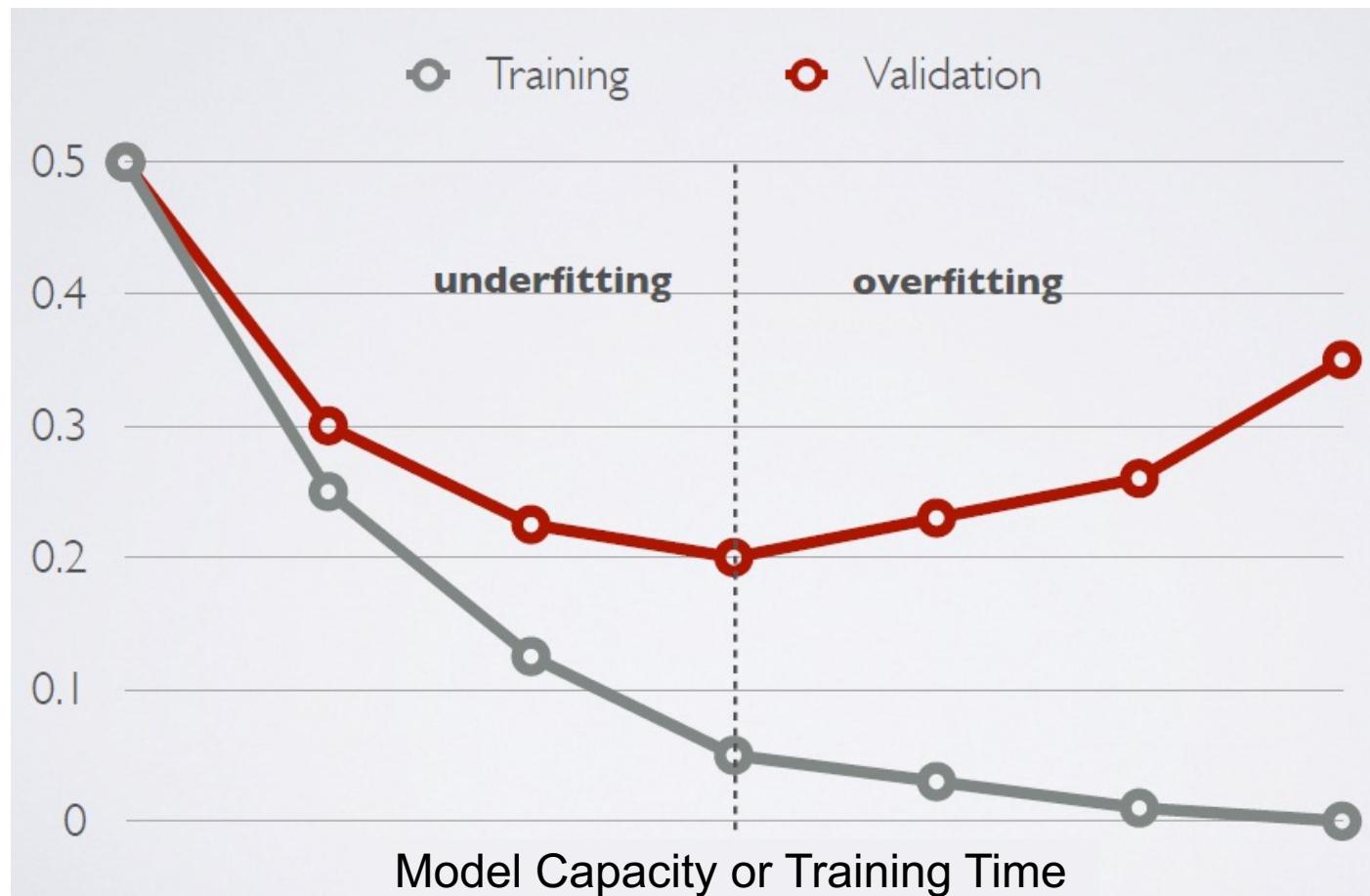
- These 3 subsets must be **perfectly disjoint** and all **representative** of the data.
- To achieve this, the split is done **at random**.

## Overfitting



- This separation is **absolutely essential** for any supervised machine learning algorithm to reliably work
- The model parameters are only optimized over the **training set**
- We use the validation set to:
  - Verify that the model is making progress **while training it**
  - Tune **hyperparameters**
  - Compare different families of models (**model selection**)
- Looking at the test set if **forbidden** in any of those steps

## Overfitting vs. Underfitting



## Example of technique to reduce overfitting

→ Data augmentation



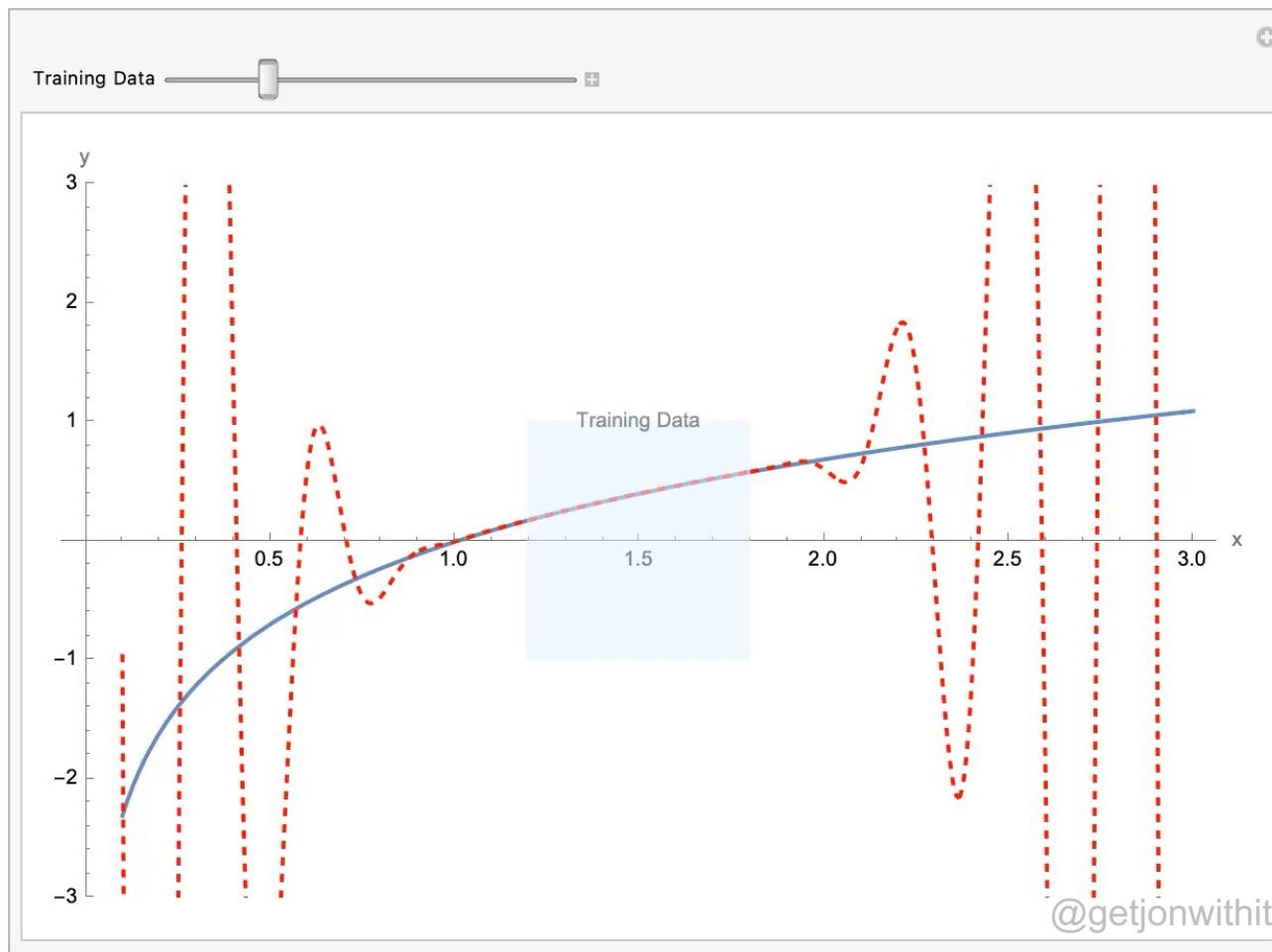
Signal transformations / Noise



Simulate data

**Exercise:** *Think of data augmentation strategies for your favorite acoustic application!*

- Machine learning models are **high-dimensional curve fitters**



- Don't expect them to **extrapolate** much beyond your training data

## The « Garbage in Garbage Out » principle

→ A machine learning model will only be as good as the data you give it

### Beware of data biases

#### Representation bias



#### Measurement / Information bias



« Please label these songs according to their genre »

#### Exclusion bias



#### Prejudice bias



#### Automation bias



#### Historical bias



Impact of smartphone usage on school grades

- Machine learning aims at finding **models** from **data** to achieve **tasks** via **optimization** over **parameterized families** of models
- Comes in **flavours**: supervised, unsupervised, predictive, generative, depending on the type of available data and task
- DNNs form a **versatile class of parameterized families of nonlinear functions** that can extract complex features from data
- **Inspired** (but far from matching!) biological brains
- DNNs are trained using variants of (stochastic) **gradient descent** and the **backpropagation algorithm**
- Made easy by modern deep learning framework such as **PyTorch**
- **Machine learning is only as good as the data you have !**

*Now let's code !*



**A S S A**  
AUTUMN SCHOOL SERIES  
IN  
ACOUSTICS